

Hypermnnesia: Eventual Consistency in Mnesia

trovato@corporation.com
trovato@corporation.com

trovato@corporation.com
trovato@corporation.com

ABSTRACT

Mnesia is a soft real-time embedded Database Management System written for Erlang, a programming language that powers the infrastructures of various organisations like Cisco, Ericsson and the NHS. Due to Mnesia’s tight integration with Erlang, it is also impactful in open source projects such as RabbitMQ and ejabberd.

However, the development of Mnesia has remained stagnant for years, resulting in the lack of features such as automatic conflict resolution: Mnesia leaves the handling of conflicts after network partitions entirely to the developer. Moreover, as a distributed database, Mnesia only provides two extreme forms of consistency guarantee: transactions and weak consistency. Existing solutions to this problem are either external libraries or commercial standalone products, none of which is integrated into Mnesia natively. This means Erlang developers often have to introduce new dependencies into their codebase or resort to less ideal alternative databases.

To address this issue, we propose a new consistency guarantee for Mnesia: eventual consistency (EC). The benefit of this is twofold: first, EC introduces an intermediate consistency guarantee between transactions and weak consistency, offering more choices to developers; second, this implementation of EC with Conflict-free Replicated Data Types (CRDTs) enables automatic conflict resolution after a network partition.

We have implemented EC as an extension to Mnesia named *Hypermnnesia* and evaluated its correctness, efficiency and usability. Evaluation results show that Hypermnnesia’s EC operations can produce consistent results in the presence of partitions and perform more than 10 times faster than Mnesia’s default transactions. Moreover, Hypermnnesia’s API enables minimum code refactoring for adoption in real-world systems. We hope Hypermnnesia can be integrated into Mnesia and used by Erlang developers in the future.

PVLDB Reference Format:

. Hypermnnesia: Eventual Consistency in Mnesia. PVLDB, 14(1): XXX-XXX, 2020.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at URL_TO_YOUR_ARTIFACTS.

1 INTRODUCTION

This project report presents *Hypermnnesia*, an extension to the Mnesia Database Management System (DBMS) [19]. Hypermnnesia provides eventual consistency (EC) guarantee on top of Mnesia’s transactional semantics. In this chapter, I first discuss the motivation for building Hypermnnesia (§1.1), then list the challenges and contributions in building Hypermnnesia (§§1.2 and 1.3). After that, an outline of the report structure is given (§1.4).

1.1 Motivation

Mnesia is an embedded distributed Database Management System (DBMS) designed for industrial-grade communication systems written in Erlang [19]. It is influential in its usage across companies and organisations like Cisco, Goldman Sachs and Mastercard [11]. It is also used in open source highly scalable message brokers and XMPP services, such as RabbitMQ [55], MongooseIM [22] and ejabberd [43]. Furthermore, Mnesia is part of the Linux, Yaws, Mnesia, Erlang (LYME) software bundle for building dynamic web pages that can handle heavy traffic [61].

Despite its usage in heavy-load applications such as WhatsApp [57] that handle tens of thousands of messages per second [34], the development of Mnesia has stayed relatively stagnant for many years¹. Due to Mnesia’s highly specialised nature and complex codebase, extending it often requires a combination of research knowledge and language expertise. This has deterred many developers from contributing to Mnesia, which means Mnesia lacks features that many modern distributed databases have. Examples include the lack of non-transactional guarantees and automatic recovery after network partitions.

Mnesia provides two ways to access the database: transactions and dirty operations. Transactions provide ACID guarantees [26], while dirty operations only give weak consistency [56]. Developers have to choose either the strong transactional API while sacrificing performance and availability during network partitions, or use dirty operations for fast performance, but risk their data being inconsistent. Moreover, Mnesia leaves the handling of potential data inconsistency after a partition entirely up to the developer [19], which means a manual restart of the cluster is often needed. Indeed, developers have reported that “*the experience has convinced us that we need to prioritize up our network partition recovery strategy*” [40] while handling a partitioned Mnesia cluster.

There are plenty of consistency models that are weaker than transactions but stronger than dirty operations’ weak consistency [18]. One such model is eventual consistency [56], which guarantees convergence eventually when no updates are made to the data. This is an attractive model for Mnesia since it often has better performance than transactions, while also maintaining consistency eventually. Modelling data this way can serve as a viable approach for achieving data consistency after a network partition.

¹Much of the core code modified in this project are over 13 years old, dating back to the time when Erlang was published on Github.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

To this end, I propose the following research questions (RQs):

- RQ1. *Can we introduce a non-transactional consistency guarantee, for example, eventual consistency, into Mnesia?* Mnesia was developed and designed before many of these consistency models became popular [37], and therefore retrofitting them into Mnesia can be a non-trivial task.
- RQ2. *How much overhead, in terms of time and space, will this new guarantee bring, compared to Mnesia's existing operations?* Mnesia's tight integration with Erlang gives it outstanding performance compared to other standalone databases. Therefore, it is important to evaluate the overhead of the new operations to ensure they are still competitive.
- RQ3. *Can eventual consistency help us achieve automatic conflict resolution after a partition in Mnesia?* The fault tolerance model of Mnesia (§4.10) may present extra challenges in this otherwise straightforward extension.
- RQ4. *How much code refactoring is needed for the consistency guarantee to function?* Mnesia is used in many large-scale (legacy) codebases, which might be too expensive to do refactoring. Therefore it is important to understand the cost of adopting a new consistency model API in Mnesia.

1.2 Challenges

There are several technical challenges (TCs) that need to be overcome to design, implement and evaluate Hypermnesia:

- TC1. Mnesia is a rather old database system that has not received much attention for years. This means the codebase contains legacy code and is not well documented.
- TC2. The previous point implies that there are potentially legacy projects using Mnesia. The new feature, therefore, needs to be backwards compatible and ideally easy to be adopted by developers.
- TC3. Another consequence of item TC1 is that its test suite does not include network partition tests, which are essential in testing the correctness of eventual consistency. Moreover, its benchmark suite is only written for transactions. Both of these need to be extended to evaluate Hypermnesia.
- TC4. There are multiple ways to achieve eventual consistency. Therefore an informed choice needs to be made based on Mnesia's existing architecture. A survey of existing methods and a study of the relatively undocumented architecture of Mnesia is essential.

1.3 Contributions

- **A new eventual consistency model for the Mnesia DBMS.** A new set of APIs is designed for this new model to minimise the amount of code refactoring for existing codebases. Moreover, with this new API, automatic conflict resolution is possible after a network partition, relieving developers from the burden of manually restarting the database.
- **An extension of the test suite for network partitions.** These are used to test the correctness of Hypermnesia. It might be useful for testing existing functionalities such as transactions as well.

- **An extension of the existing benchmarking library for dirty operations and EC operations.** These are built on top of the original benchmarks for transactions.
- **An evaluation of the new eventual consistency API,** in terms of the following: (a) fault tolerance by testing its correctness under network partitions and ability to reconcile after a partition; (b) performance by analysing time and space overhead and comparing against existing operations. Results show that Hypermnesia's EC operations can often achieve 10x higher throughput and lower latency than Mnesia's transactions; (c) usability by analysing the amount of code refactoring needed to adopt the new API in real codebases.

1.4 Outline

The rest of the report is organised as follows: §2 introduces the reader to the background knowledge, including the architecture of Mnesia and different ways of achieving eventual consistency. ?? describes recent research on eventually consistent databases and CRDTs: a data structure for implementing eventually consistent systems. §4 goes into detail about designing an API for Hypermnesia and implementing it with practical optimisations. §6 evaluates Hypermnesia in terms of its correctness, performance and usability. §7 concludes the report and outlines potential future work.

2 BACKGROUND

This chapter introduces consistency models in a distributed system relevant to this work (§2.1). It then discusses the architecture of Mnesia (§2.2) before giving an overview of CRDTs.

2.1 Consistency models

A distributed system often involves a set of replicated state machines [33]. Coordinating these clusters of machines is a challenging task, and tradeoffs can be made between consistency and performance. Many consistency models are formalised for both transactional and non-transactional systems [54]. This section gives a brief introduction to three of them in detail: distributed transactions (provided by Mnesia's transaction API), weak consistency (provided by Mnesia's dirty operations) and eventual consistency (the aim of this project).

Distributed transactions. is sometimes also called distributed atomic commits. It is achieved if all nodes commit or all nodes abort [47]. This is part of the properties provided by ACID [26] and is often implemented with the distributed two-phase commit protocol [9]. Mnesia provides such a guarantee via its transaction APIs.

Weak consistency. is defined as follows: "The system does not guarantee that subsequent accesses will return the updated value, and several conditions need to be met before the value will be returned" [8, 54, 56]. It is (deliberately) vaguely defined to incorporate systems whose replicas "might become consistent by chance" [8]. Mnesia's dirty operations fall into this category.

Eventual consistency. is defined as follows: "If no new updates are made to the object, eventually all accesses will return the last

updated value” [56]. There are often two steps in achieving eventual consistency [62]: (i) anti-entropy [12] disseminates data to all the nodes in a cluster; (ii) conflict resolution addresses potential conflicts when handling received data.

Anti-entropy typically involves some form of broadcasting messages, while conflict resolution tends to use (one of) the following three techniques: [29]: *Conflict-free Replicated Data Types (CRDTs)* [41, 50] *Mergable persistent data structures* [23] and *Operational transformation* [13].

Among these three ideas, CRDTs is one of the database community’s most used and well studied approaches. It has been successfully deployed in many NoSQL databases, such as Riak [32], Redis [46] and Microsoft Azure Cosmos DB [51]. Operational transformation requires a central coordinator, which is unsuitable for Mnesia as it uses a leaderless replication strategy (§2.4). Mergable persistent data structures are based on version control ideas, but multi-version support is not available in Mnesia, meaning this method is not easily implementable in Mnesia. In conclusion, CRDTs is chosen as the basis of conflict resolution and discussed in more detail in §2.8.

2.2 Mnesia

Mnesia is an embedded [44], specialised distributed DBMS for Erlang/OTP applications (§§ 2.3 and 2.4). It is similar to relational databases where each table stores tuples (§2.5), but also different since it uses Erlang as its query language instead of SQL. It has dual mode support for accessing data (§2.6), discussed in more detail below.

2.3 Design goals

The original design of Mnesia attempts to meet several requirements [19]:

- (1) Fast real-time key-value lookup
- (2) Complex non-real-time queries (mainly for operation and maintenance tasks)
- (3) Distributed data (due to the distributed nature of the applications)
- (4) High fault tolerance
- (5) Dynamic reconfiguration

Based on these, Mnesia was built as part of the Erlang/OTP distribution. Erlang was originally designed for building fault-tolerant telecommunication systems, and Mnesia helps it to better achieve that goal by acting as an embedded database: it is tightly coupled to Erlang, giving it two distinct features: (a) The database runs in the same address space as the application itself, offering minimum overhead while accessing data. (b) The database uses native Erlang records to store its data, removing the impedance mismatch between different data formats.

These special features make Mnesia only suitable for specific purposes. Mnesia is typically used when there is a need to replicate a relatively small amount of data: compared with standard SQL databases that can handle terabytes of data, Mnesia is instead built for (tens of) gigabytes of data [27]. For example, user login details are often stored as session data, and to scale the application out, these data need to be replicated across nodes and accessed quickly to

provide a good user experience. Mnesia can be a suitable candidate in such a case due to its compelling performance.

2.4 Architecture

Mnesia is built on top of Erlang’s built-in memory and disk term storage ets and dets [21]. These term storage can be thought of as primitive storage engines that provide constant (or logarithmic) access time for large amounts of data [27]. They support different data structures for storing data, such as set, bag, etc. Internally, these are implemented as hash tables or balanced binary trees. Mnesia also provides additional functionalities such as transactions and distribution on top of ets and dets.

A Mnesia cluster generally has a leaderless architecture where every replica can handle client requests. A cluster of Mnesia nodes are connected via the Erlang distribution protocol, which uses TCP/IP as its carrier by default, providing reliable in-order delivery. Moreover, the connection is transitive, like fig. 1 which forms a cluster of fully connected nodes (or a mesh).

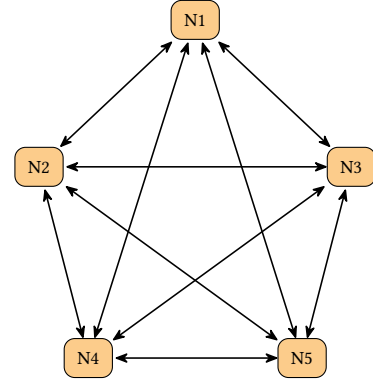


Figure 1: Example Mnesia cluster of five nodes. Note they always form a fully connected network.

Figure 2 shows a brief overview of the modules in Mnesia according to the supervision relation. Erlang’s fault tolerance model has the concept of a supervision tree [20], a hierarchy of processes in charge of monitoring/supervising child processes for killing and restarting “misbehaving” processes.

2.5 Data representation

Mnesia stores data in Erlang records [21]. Table 1 is an example of a student record, and this is how each row of a Mnesia table is represented, with listing 1 as the corresponding definition in Erlang.

student	id (key)	name	college	age
	bb123	Bruce Banner	Avengers	54
	ts233	Tony Stark	Avengers	50
	sg333	Steve Rogers	Avengers	100

Table 1: An example Mnesia table for student information.

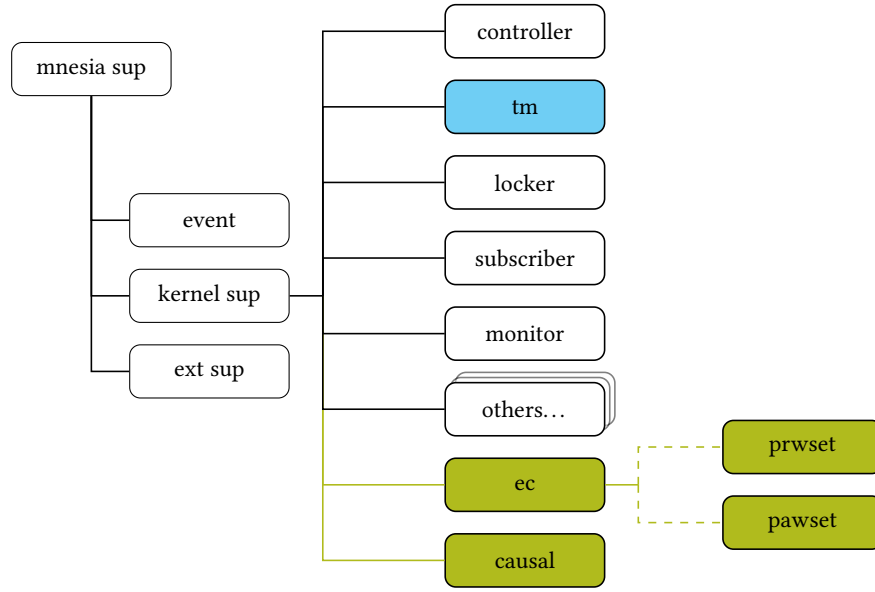


Figure 2: Relationships between various Mnesia processes, in the form of a supervision tree. Most of the logic for replication is inside `mnesia_tm` (highlighted in blue). Modules added by Hypermnesia are highlighted in green. More details on the new modules are discussed in §4.1. Diagram partly based on [38].

```
-record(student,
  {id :: integer(),
   name :: string(),
   college :: string(),
   age :: integer()}).
```

Listing 1: Erlang record definition for a student.

2.6 Access contexts

A central API provided by Mnesia for table manipulation is given in listing 2. A user calls the `activity` function which takes in:

- (1) A `Kind` of activity, also called a *access context*. Currently supported contexts include transactions and dirty operations;
- (2) A function (`Fun`);
- (3) Arguments to be passed to the function (`Args`);
- (4) The module in which the function is defined (`Mod`).

There are two almost equivalent² ways of using this API, and listing 3 provides an example of writing a tuple `{k,v}` into the table `tab_name` and then reading it out. I use the convention of `function_name/arity` to represent Erlang functions in these code examples and the rest of this report.

```
mnesia:activity(Kind, Fun, Args :: [Arg :: term()], Mod) ->
  <- t_result(Res) | Res
  where Kind = activity()
        Fun = fun(...) -> Res
        Mod = atom()
```

Listing 2: Mnesia table manipulation API.

²Subtle differences between these two APIs can be found in the reference guide [18]

```
mnesia:activity(transaction,
  <- fun () ->
    mnesia:write({tab_name, k,
                  <- v}),
    mnesia:read({tab_name, k})
  end).
```

(a) transaction with activity/2

```
mnesia:transaction(fun () ->
  mnesia:write({tab_name, k,
                <- v}),
  mnesia:read({tab_name, k})
end).
```

(c) transaction with transaction/1

```
mnesia:activity(async_dirty,
  <- fun () ->
    mnesia:write({tab_name, k,
                  <- v}),
    mnesia:read({tab_name, k})
  end).
```

(b) asynchronous dirty with activity/2

```
mnesia:async_dirty(fun () ->
  mnesia:write({tab_name, k,
                <- v}),
  mnesia:read({tab_name, k})
end).
```

(d) asynchronous dirty with async_dirty/1

Listing 3: Comparing transaction and asynchronous dirty operations API.

Transactions. in Mnesia provide ACID properties [19]. While this is attractive, it also incurs large performance overhead and hence a full-fledged transaction might not be suitable for all tasks. For example, in a datagram routing application, it can be too slow to initiate a transaction each time a packet is received [19]. The underlying implementation of Mnesia uses two-phase commit for distributed transactional atomicity, write-ahead logging for durability and consistency, and two-phase locking for isolation.

Dirty operations. in Mnesia bypass most transactional processing and operate directly on the data. Therefore they are much faster (at least 10x [18]) than transactions. As a consequence, they lose the

ACID properties. However, dirty operations still provide some level of consistency such as no garbled records for individual dirty reads. The underlying implementation of dirty operations follows this pattern: when the function is called, the operation is first performed on the local table. Then relevant information is collected including the Erlang record, the target table, and so on. These are then sent to other replicas in the cluster. Mnesia does not examine the content in the table during this process. This property is useful in helping us choose a CRDT, as we will see later that some CRDTs require examinations of the state before broadcasting the message (§3.2). On the receiving side, the transaction manager monitors and applies the received message accordingly.

2.7 Consistency and availability

The CAP theorem [24] forces a system to choose between consistency and availability during a network partition. Mnesia responds to this problem by taking *no* stance with respect to the CAP theorem [2], but leaves the handling of the recovery process entirely to the developer [19], often resulting in a manual restart of the Mnesia nodes after partitions. Recent improvements have incorporated the feature to allow developers to set the `majority` option which disallows any non-dirty operation to a table not in a majority quorum. This option is useful for preventing conflicts for critical data, but sacrifices the availability of the database further and only works for transactions. There is also a `set_master_nodes/2` function which unconditionally loads data from the master node after a partition. It can be useful in some cases, but is a rather “brutal” way of handling partitions since it might lead to data loss.

Erlang’s built-in distribution protocol has a failure detector that periodically sends heartbeats. If a node does not receive a heartbeat from another node for some time, then the failure detector considers that node dead and disconnects itself. Mnesia relies on Erlang’s failure detector to function, and two possible situations can occur during a network partition (NP):

- NP1. The partition is a *transient failure* in which the failure detector does not consider a temporarily unresponsive node dead. In this case, operations can carry on as usual, but transactions will stall since the two-phase commit protocol requires responses from all participating nodes. Dirty operations can be performed as usual. When the partition heals, buffered messages will then be delivered. But dirty operations risk the database being in an inconsistent state when the network heals, as we shall see in §4.11.
- NP2. The partition lasts long enough that the failure detector detects the *communication failure* and disconnects the nodes it thinks have failed. In this case, Mnesia would reconfigure itself and perform future operations without the failed nodes, i.e. Mnesia considers itself as having a new cluster with only the members still alive. When the partition heals, Mnesia logs an error message saying that the database might be inconsistent and ask the developer to fix this issue manually.

2.8 CRDTs

Conflict-free Replicated Data Types (CRDTs) [41, 50] are a family of replicated data types with a common set of properties that enable

operations to be performed locally on each node while always converging to a final state among replicas if they receive the same set of updates. There are two types of CRDTs: state based (§2.9) and operation based (§2.10). Each of them is discussed in detail with examples below.

2.9 State-based CRDTs

Here I provide a summary of the properties of state-based CRDTs. Refer to ?? and [1, 41, 50, 53] for mathematical details and examples.

A state-based CRDT communicates by propagating its state of the data structure with other parties, and performs the merge operation with a merge (sometimes called join) operator to converge towards the Least Upper Bound (LUB) of two states. The merge operator must be commutative, idempotent and associative to guarantee the convergence property (proposition 2.1), which is key to eventual consistency.

PROPOSITION 2.1 ([50]). *Any two object replicas of a state-based CRDT eventually converge, assuming the system transmits payload infinitely often between pairs of replicas over eventually-reliable point-to-point channels.*

Propagating the entire states of a (big) data structure can often be expensive. Therefore Delta-State Conflict-Free Replicated Data Types (δ -CRDTs) [1] is proposed. It disseminates only the changing parts of the state as a δ -state, reducing communication costs. Due to their simple requirements on the channel, δ -CRDTs often have complex logic in managing and disseminating states.

2.10 Operation-based CRDTs

Operation-based (op-based) CRDTs send the operation performed on the CRDT object (as opposed to the state) along with some metadata. An op-based CRDT typically requires two *concurrent* operations f and g to be commutative under the causal delivery order $<_d$: $f \parallel_d g \iff f \star_d g \wedge g \star_d f$.

For this reason, a causal delivery channel is often required for op-based CRDTs to work. Given such a channel, the following property (proposition 2.2) holds:

PROPOSITION 2.2 ([50]). *Any two replicas of an op-based CRDT eventually converge under reliable broadcast channels that deliver operations in delivery order $<_d$.*

One advantage of op-based CRDTs is their communication efficiency. Using a set as an example data structure, instead of sending the entire set with all its elements, op-based set can just send operations like add or delete. This could significantly reduce the communication cost, especially when the set is large.

Given the stronger assumption on the channel, designers of op-based CRDTs only need to choose how they want to order concurrent operations, an example of a particular type (which is also the one implemented in Hypermnesia) of op-based CRDTs is given below.

2.10.1 Pure op-based CRDTs. The definition of the original op-based CRDTs makes the separation between state-based and op-based CRDTs less clear in that an op-based CRDTs can often include state information while sending operations. For example, before

sending operation `add(1)` on a set, the state of the set can be inspected and the entire set can be added as the “operation”. It also has a relatively high implementation complexity [6]. To address this, pure op-based CRDTs were developed [4, 5], which forbids inspecting the state σ while generating the message m .

I now illustrate pure op-based CRDTs with an example add-wins set (AW-set). Listing 4 is the conceptual view of implementing it. Typically there are three operations on an op-based CRDT: a generator `prepare`, which turns the operation o_i on the current state σ_i into a message m_i that is sent via reliable causal broadcast; an effector `effect`, which receives the message m and applies it onto the current state σ_j and produces a new state σ'_j ; and `eval`, which takes an operation, such as lookup an element in a set and the current state σ_i , and returns the result of the operation. For pure op-based CRDTs, data and timestamps need to be stored in a *Partially Ordered Log (PO-Log)* [4, 5], partially ordered by the timestamp attached to each operation. This aims to identify causal relationships between different elements, which is important for the convergence property (proposition 2.2).

```

 $\Sigma : T \hookrightarrow O$ 
preparei( $o, s$ ) =  $o$ 
effecti( $o, t, s$ ) =  $s \cup \{(t, o)\}$ 
evali(elems,  $s$ ) =  $\{v \mid (t, [\text{add}, v] \in s \wedge \nexists (t', [\text{rmv}, v] \in s \cdot t < t'))\}$ 

```

Listing 4: A pure op-based AW-set implementation, adapted from [5].

The PO-Log (Σ) for the AW-set is a partial function from timestamps to operations (with appropriate payload). Note that the effector takes three arguments, the operation o , timestamp t and state s , and it simply unions the operation and timestamp with the current state (i.e. appending it to the log), leaving the complexity to the lookup function. This is a fairly space-expensive operation, since all operations need to be stored in the log including deletions. Optimisations are discussed later in §4.5.

2.11 Summary

In this chapter, we discussed transactions, weak and eventual consistency (§2.1), the first two are offered by Mnesia (§2.2), and this project aims to provide the last. A family of data structures, (CRDTs) are also described (§2.8) to show how they help achieve eventual consistency.

3 SYSTEM DESIGN

This chapter provides an account of Hypermnnesia’s design and implementation. Hypermnnesia’s API and the selection of a CRDT are discussed in ??, with alternative options also considered. Subsequently, the implementation is outlined in §4.1 at a high level before diving into three key components: causal broadcast, Set CRDT and fault tolerance from §§4.2 to 4.10.

3.1 API design

Hypermnnesia’s API is designed with the research question item RQ4 on refactoring in mind. We consider the SOLID principle [36] with the following design goals (DGs):

- DG1. **Backwards compatibility and code refactoring.** Hypermnnesia should be backwards compatible with the existing codebase and should run normally without modification if developers choose not to use the new feature. Moreover, if they opt for the new feature, the API should minimise the amount of refactoring needed for an existing codebase.
- DG2. **Single responsibility.** The new feature should be contained in module(s) isolated from existing Mnesia code, while respecting Mnesia’s code structure. This makes the implementation self-contained and easier to maintain.
- DG3. **Extensibility.** Hypermnnesia should be extensible to new implementations. This is beneficial since there are various possible CRDTs, e.g. operation-based and state-based. They often achieve similar goals but rely on different assumptions. Hypermnnesia should be extensible to new CRDT variants for different tradeoffs.

Based on the existing APIs provided by Mnesia (§2.6), one way to extend it with new APIs for eventual consistency would be to add a new access context called `async_ec` so that database operations performed within this context are eventually consistent. Using the same example, we can change listing 3 (the original API for transactions and dirty operations) to listing 5 (the new EC API). We argue that this is a fairly natural extension to the existing Mnesia access contexts, and requires little refactoring of existing code (dirty or transaction changed to `ec`). It does not break existing features either if no changes are made to the existing code, thus fulfilling item DG1.

<pre> mnesia:activity(async_ec, fun () ↪ -> mnesia:write({tab_name, k, ↪ v}), mnesia:read({tab_name, k}) end). </pre>	<pre> mnesia:async_ec(fun () -> mnesia:write({tab_name, k, ↪ v}), mnesia:read({tab_name, k}) end). </pre>
<p>(a) <code>activity/2</code> with new access context <code>async_ec</code></p>	<p>(b) <code>async_ec</code> with new function <code>async_ec/1</code></p>

Listing 5: New eventual consistency (EC) API based on existing Mnesia APIs, using the same example as listing 3.

We also wish to allow programmers to declare which type of CRDTs they want to use. This can be done while declaring the table type as `pawset` at creation time. Mnesia asks users to enter the type of each table at creation time (?? 6a). The default values are set, bag, etc. This can be extended to support pure AW-set (`pawset`) and RW-set (`prwset`) as well (?? 6b). Note that having multiple Set CRDT implementations allow developers to choose the most appropriate one for their applications, demonstrating the extensibility of Hypermnnesia (item DG3).

Finally, different CRDT logic is implemented inside separate modules, aiming to achieve the single responsibility requirement (item DG2). The new `mnesia_ec` module is built for handling data


```

mnesia:create_table(project,
  [{type, set},
   {ram_copies,
    ↪ all_nodes()},
   {attributes,
    ↪ record_info(fields,
    ↪ student)}}]).

mnesia:create_table(project,
  [{type, pawset},
   {ram_copies,
    ↪ all_nodes()},
   {attributes,
    ↪ record_info(fields,
    ↪ student)}}]).

```

(a) Creating a table of default type set and storing student record data listing 1.

(b) Creating a table of type pawset (pure add-wins set) and storing student record data. A prwset (pure remove-wins set) can be used as well.

Listing 6: New eventual consistency (EC) API based on existing Mnesia APIs.

replication and interfacing with the underlying CRDTs, `mnesia_causal` module for causal delivery, and `mnesia_pawset/mnesia_prwset` module(s) for data storage and conflict resolution (fig. 2).

3.2 Choosing a CRDT

We talked about different kinds of CRDTs in §2.8. One type of Set CRDT needs to be chosen for Hypermnesia’s implementation. We consider two main factors (Fac):

- Fac1. How efficient is this CRDT in terms of space and time?
- Fac2. How easy does it fit into the existing Mnesia codebase and how many breaking changes need to be made?

On the one hand, state-based CRDTs (§2.9) have an important drawback in their communication overhead [1]. This might not be acceptable for data types with large state such as sets. δ -CRDTs is a more suitable candidate for our purpose but even with δ -CRDTs, a fair amount of data needs to be broadcast, especially as the number of operations increases (see ?? for details). Moreover, it might be difficult to examine the state before broadcasting, since Mnesia does not do this by default (§2.6). State-based CRDTs do have the advantage of low demand on the channel. Therefore reliable broadcasting in Mnesia would be sufficient.

On the other hand, operation-based CRDTs tend to require the channel to provide causal broadcast, which has to be implemented in Mnesia from scratch. Dynamic membership, i.e. nodes leaving and joining the cluster, is also trickier with operation-based CRDTs as buffering and replaying of operations are needed. Nevertheless, pure operation-based CRDTs largely resemble Mnesia’s anti-entropy strategy, with no examination of the current content and immediate synchronisation for each operation [39]. State-based CRDTs are less suitable for these requirements [41]. Therefore pure op-based CRDTs better meets the second requirement (item Fac2).

In terms of efficiency, it is challenging to characterise the performance metrics of each CRDT without practical benchmarking: δ -CRDTs uses periodic synchronisation of delta states, while pure op-based CRDTs go through an extra causal broadcast layer. Moreover, Almeida et al. [1] and Baquero et al. [5] proposed optimisations for these CRDTs. One could also argue that these two CRDTs are fundamentally the same as they are both doing the necessary work for conflict resolution, but at different layers of the system.

In conclusion, pure op-based CRDTs is chosen for its operational similarity with Mnesia’s dirty operations. We leave experimenting of δ -CRDTs as future work (§7.2). Table 2 highlights the differences between these two CRDTs.

	State-based
Variant of interest	Delta-state
Understandability	Complex
Efficiency	Good
Complexity	Periodic anti-entropy different from Mnesia’s protocol
Advantage	Little requirement on the channel

Table 2: Comparison of state and operation based Set CRDTs.

4 SYSTEM IMPLEMENTATION

4.1 Implementation overview

Figure 3 shows an overview of Hypermnesia’s architecture. A client sends a request to one of the database nodes connected in the cluster, and this request is then processed by the `mnesia_ec` module (§§2.4 and 3.1), responsible for calling the underlying pure op-based Set CRDT implementation. In Hypermnesia, two op-based Set CRDTs are implemented: add-wins set (`mnesia_pawset`) and remove-wins set (`mnesia_prwset`) (§4.5), which are called by `mnesia_ec` based on the appropriate table type (listing 6). `mnesia_ec` also calls the `mnesia_causal` module that handles the causal broadcast (§4.2). Each of these modules is explained in more detail in the following sections.

4.2 Causal Broadcast

Causal broadcast is a mature algorithm with standard implementation strategies [10, 48]. This section focuses on various modifications of the causal broadcast algorithm and pure op-based set to the Erlang ecosystem.

4.3 Causal broadcast server

Causal broadcast is often treated as a middleware between the network and the application. It buffers messages until they are causally ready to be delivered to the application. The `gen_server` behaviour is a suitable abstraction, as it provides an interface from which developers can write custom implementations of request handlers to obtain a generic server (similar to interface in object-oriented languages) [20].

More concretely, we define a collection of public functions for the causal broadcast server, which can be called by `mnesia_ec` while sending and receiving messages. Two most basic ones are `send_msg/1` and `rcv_msg/1`, as shown in listing 7. The `send_msg/1` function returns the current timestamp as a vector clock to be attached to the message. The `rcv_msg/1` function takes a received message and returns a list of messages ready to be delivered to Mnesia for further processing, such as database writing. The `mnesia_causal` server also

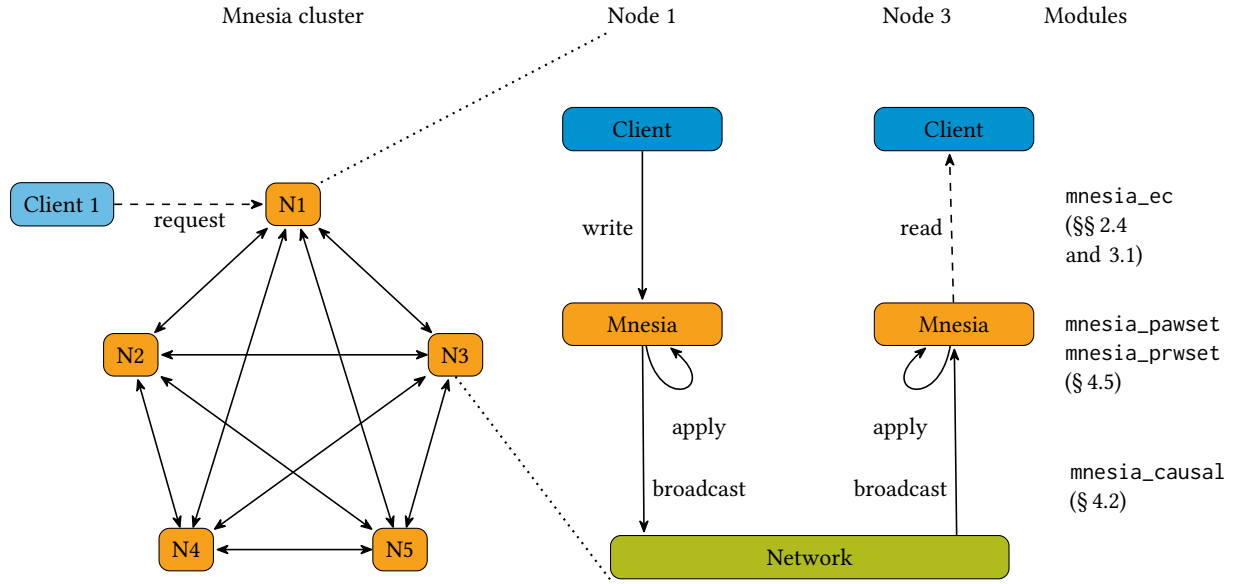


Figure 3: Hypermnesia architecture overview.

keeps track of a list of buffered messages. Each time a message is received, it is added to the buffer, and the buffer is searched for deliverable messages.

```
-record(state,
  {send_seq :: integer(),
   delivered :: vclock(),
   buffer :: [msg()] }).

-spec send_msg() -> vclock().
-spec rcv_msg(msg()) -> [msg()].
```

Listing 7: mnesia_causal server for causal broadcast.

4.4 Tagged causal stable broadcast

The algorithm discussed above is a standard causal broadcast algorithm [10, 30, 48]. It does not expose any ordering or timestamp information when delivering a message. A Set CRDT often relies on unique identifiers to ensure commutativity of causally concurrent operations. Baquero et al. [5] observe that the timestamp information from the causal broadcast layer can act as a unique identifier for the CRDT, and thus propose exposing the ordering information from the causal broadcast layer to the application, i.e. the receive function would return a list of messages along with their timestamps (listing 8).

```
-spec send_msg() -> timestamp().
-spec rcv_msg(msg()) -> [{msg(), timestamp()}].
```

Listing 8: mnesia_causal server interface for Tagged Causal Stable Broadcast.

The extended Tagged Causal Stable Broadcast (TCSB) protocol works like this: when a message is ready to be sent, the function `send_msg/0` is called to obtain the timestamp which is attached to the message. And when a message is received, the `rcv_msg/1` is called to buffer the message and find all messages ready to be delivered. A message is ready when the following two conditions hold [10, 30, 48]:

- (1) The sender entry in the receiver's delivered map, which keeps track of how many messages are delivered for each sender, is exactly one less than the entry in the received message's timestamp;
- (2) The other entries in the received message's timestamp are all less than or equal to the corresponding entries in the receiver's delivered map.

The first condition ensures that all causally preceding messages from the sender have been delivered, and the second condition prevents us from delivering messages that causally depend on messages from other nodes. Note that the sender of the message can always deliver messages from itself immediately, a desirable property for availability during partitions, and causal consistency is the strongest consistency model that provides such always-available property [35].

4.5 Pure operation-based Set CRDT

We have discussed the idea of a pure op-based set (§ 2.10). This section builds on top of the basic ideas of op-based Set CRDT and highlights a few optimisations exploiting the causal broadcast layer (§§ 4.6 to 4.8). More subtle implementation challenges in interactions between Mnesia and the pure op-based set CRDT are also discussed (§ 4.9).

4.6 Causal Redundancy

Causal redundancy is proposed by Baquero et al. [5] for elements in the PO-Log that can be removed without affecting the output of queries. We define the redundancy relation in listing 9. The first rule says that a clear or remove operation is redundant in all cases. The second rule makes causally older additions of the same elements redundant. The third rule says every addition of an element present in the PO-Log is in the set. Redundant elements can be safely removed from the PO-Log, which helps reduce the storage cost.

$$\begin{aligned} (t, o) R s &\iff o[0] = \text{clear} \vee \text{rmv} \\ (t, o) R_-(t', o') &\iff t < t' \wedge (o'[0] = \text{clear} \vee o'[1] = o'[1]) \\ \text{eval}_i(\text{elems}, s) &= \{v \mid (_, [\text{add}, v]) \in s\} \end{aligned}$$

Listing 9: The redundancy relation R for an AW-set. $o[0]$ is the operation, while $o[1]$ is the key. Taken from [5].

4.7 Causal stability

The number of timestamps associated with each element grows as elements are added to the set. Moreover, each of them is proportional to the number of nodes in the cluster. Baquero et al. [4] suggest removing the timestamps that are *causally stable* [5, 50]:

Definition 4.1 (Causal Stability [4, 5]). A timestamp τ and a corresponding message, is causally stable at node i when all messages subsequently delivered at i will have timestamp $t > \tau$. Or equivalently:

$$\text{tcstable}_i(\tau) \text{ if } \forall j \in I \setminus \{i\}. \exists t \in \text{deliv}_i() \cdot \text{src}(t) = j \wedge \tau < t.$$

Once a timestamp is stable, it can be removed. This is safe since there will be no more concurrent operations delivered in the future and the timestamp metadata is used to ensure commutativity for concurrent operations and is therefore no longer needed, reducing the storage cost of the PO-Log.

The actual implementation in Hypermnnesia is done by asking the pure AW-set periodically scanning elements and removing causally stable timestamps.

4.8 Broadcast and CRDT algorithms

Algorithm 1 shows the broadcast algorithm obtained by putting these optimisations together (§§4.4 to 4.7). When broadcasting a message (?? 1–1), the appropriate timestamp and sender id are attached to the message, which is then received and buffered (?? 1–1). The receiver searches for messages that are causally ready to be delivered using the condition in §4.4 (?? 1–1). The causal broadcast server also provides a function to check for the stability of a timestamp (?? 1–1). This algorithm combines the standard causal broadcast algorithm [10, 30] and extra exposed APIs proposed by [4, 5], adapted to fit Mnesia’s architecture (§2.4).

Algorithm 2 shows the corresponding algorithm for the pure AW-set, written according to the specification [4, 5]. The original specifications are written as a mathematical specification, while we

take a programming perspective and present them in terms of set operations like add, delete, etc. The `remove_redundant` function is called each time a deletion or insertion happens, and checks for redundant elements using the `redundant` function (?? 2–2), based on the idea of causal stability (listing 9).

Algorithm 1: Tagged Causal Stable Broadcast (TCSB) protocol algorithm, ideas taken from [5, 6, 30].

```

on init:
    sendSeq ← 0
    buffer ← ∅
    delivered ← {(s, 0) | ∀s ∈ nodes}
    /* a map from node to last the timestamp of
       the last delivered message */
    ts_map ← {(s, ⊥) | ∀s ∈ nodes}

on broadcast msg at replica i:
    delivered[i] ← sendSeq + 1
    sendSeq ← sendSeq + 1
    deps ← delivered
    broadcast (msg, i, deps) to all nodes

on receive (msg, sender, deps) at replica i:
    buffer ← buffer ∪ (msg, sender, deps)
    repeat
        /* find all causally deliverable messages */
        D ← {(msg, sender, τ) | ∃(msg, sender, τ) ∈
            buffer.deps[sender] = τ[sender] + 1 ∧ ∀s ∈
            dom(deps) \ {sender}. deps[s] ≥ τ[s]}
        /* update delivered vector and timestamp map */
        for (msg, sender, τ) ∈ D do
            delivered[sender] ← delivered[sender] + 1
            ts_map[sender] ← τ
        buffer ← buffer \ D
        deliverable ← deliverable ∪ D
    until D = ∅
    deliver deliverable to application

on check stability of τ:
    if τ ≤ min({ts_map[s][src(τ)] | s ∈ dom(ts_map)})
    then
        return true
    else
        return false

```

Apart from the add-wins set, a remove-wins set (§4.1) is also implemented, and the algorithm is mostly similar to that of an add-wins set. This can be found in ??.

4.9 Practical concerns

This section discusses several practical issues in implementing the algorithm.

Algorithm 2: Pure AW-set pseudocode, defined in terms of usual set operations. This pseudocode sacrifices efficiency for clarity. Ideas are taken from [4–6].

```

POLog  $\leftarrow$  []
function add( $e, t$ ):
  remove_redundant( $e, t, \text{add}$ )
  for ( $e', t', o'$ )  $\in$  POLog do
    if redundant( $((e, t, \text{add}), (e', t', o'))$ ) then
      return
  POLog  $\leftarrow$  append(POLog, ( $e, t, \text{add}$ ))

function delete( $e, t$ ):
  remove_redundant( $e, t, \text{delete}$ )
  for ( $e', t', o'$ )  $\in$  POLog do
    if redundant( $((e, t, \text{delete}), (e', t', o'))$ ) then
      return
  POLog  $\leftarrow$  append(POLog, ( $e, t, \text{delete}$ ))

function read( $k$ ):
  for ( $e, t, o$ )  $\in$  POLog do
    if  $e.\text{key} = k$  then
      return  $e$ 
  return undefined

function remove_redundant( $e, t, o$ ):
  for ( $e', t', o'$ )  $\in$  POLog do
    if redundant( $((e', t', o'), (e, t, o))$ ) then
      POLog  $\leftarrow$  remove(POLog, ( $e', t', o'$ ))

function redundant( $((e, t, o), (e', t', o'))$ ):
  /* check whether ( $e, t, o$ ) is made redundant by
     ( $e', t', o'$ ) */
  if  $o = \text{delete}$  then
    return true
  else if  $e = e' \wedge t < t'$  then
    return true
  else
    return false

periodically
  for ( $e, t, o$ )  $\in$  POLog do
    if stable( $t$ ) then
      POLog  $\leftarrow$  remove(POLog, ( $e, t, o$ ))
      POLog  $\leftarrow$  append(POLog, ( $e', \perp, \text{no-op}$ ))

end

```

Faster access to the PO-Log. As discussed in §2.4, Mnesia uses ets and dets as its storage system, which are implemented as hash tables and balanced binary trees. There is no direct support for log-like structures such as lists. This presents challenges as to how the PO-Log could be implemented. Typically a Mnesia table is a set indexed by the n -th element in a tuple, where n is customizable. We keep this structure but use a bag instead of a set to allow multiple elements with the same key but different (concurrent)

timestamps. This representation gives us the advantage of accessing the element by key in $O(1)$ time, which is useful in speeding up the remove_redundant function that needs a linear scan of the entire set for matching keys (?? 2 in alg. 2). However, this representation also loses the partial order of the PO-Log, making the stabilisation process (§4.7) much more difficult.

Payload conflicts. Algorithm 2 deals with operations of a set in terms of addition, deletion, etc. Although the Mnesia documentation claims the data structure to be a set, it behaves more like a map, with keys and payload (recall from §2.5 we saw that each row in the student table has an id as the key, and other information like name as the payload). Therefore the Set CRDT falls short when there are concurrent additions of the same key but different payloads. This can be solved in general with a Map CRDT [50], which is a straightforward extension of a Set CRDT, but requires the payload to be CRDTs as well so that the conflict can be resolved recursively for the payload.

This puts constraints on the data a user can put in a Mnesia table, and is considered as breaking too much compatibility of the existing system and therefore I take a simpler approach by resolving them based on the Erlang term order [15], an ordering of data types in Erlang. This could easily be extended with other resolution strategies as well.

More operations on Mnesia tables. Apart from the basic addition and deletion of elements, Mnesia supports a range of other operations as well, including matching based on a pattern specification, iterating over a table, finding all the keys of a table, etc. As an example, the original select/2 function in Mnesia takes a table and a pattern specification as its argument. This is modified by appending wild cards for the timestamp and operation name (since they are stored along with each tuple in the PO-Log) to the input pattern before matching. Other functions like all_keys/1 are implemented with similar ideas.

4.10 Fault tolerance

This section continues the discussion from §2.7 on Mnesia’s response to network partitions and talks about how Hypermnesia addresses them. We continue to use the term transient failure and communication failure to distinguish them.

4.11 Communication failure

Table 3 shows the sequence of operations and the corresponding states of each replica, and fig. 4 shows the graphical representation where three Mnesia nodes are initially connected to each other and each has an element a in them. Then a partition occurs between node A and node B, as well as node A and node C. During this partition, an addition of element c occurs at A while addition of element b happens at B and C. When the partition recovers, replicas are now in an inconsistent state which will be reported to the developer for manual resolution.

Hypermnesia resolves this issue by buffering messages during the partition, as shown in fig. 5c. During the communication failure, operations performed on replicas A and B are buffered until the partition heals. By the property of op-based CRDTs (proposition 2.2), replicas will converge after the buffered messages are delivered. In

Operation	Node A	Node B	Node C
A::add(a)	{a}	{a}	{a}
B::add(b)	{a}	{a,b}	{a,b}
A::add(c)	{a,c}	{a,b}	{a,b}
inconsistent_database	{a,c}	{a,b}	{a,b}

Table 3: Operations performed on each replica. Operations that happen during the network partition sit between two horizontal lines.

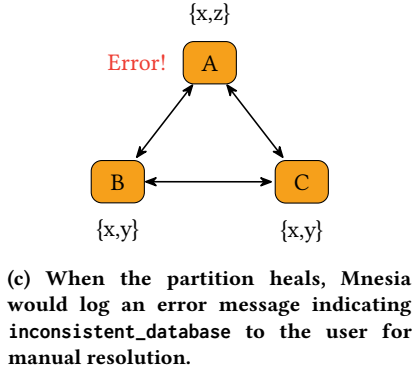
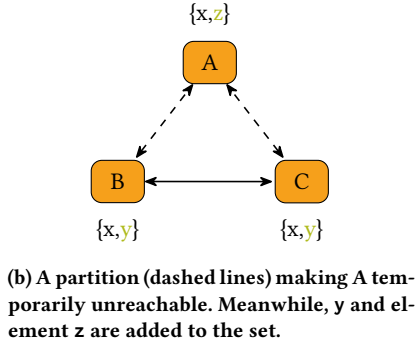
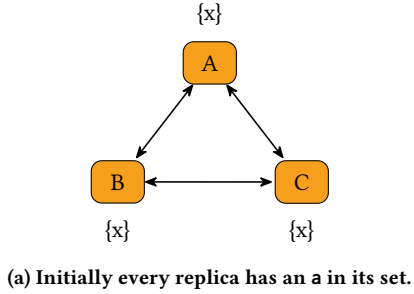


Figure 4: Illustration of a communication failure in Mnesia.

this simple case, it is sufficient to achieve consistency by a simple set union. In a more complex case such as concurrent addition and deletion, the underlying op-based CRDT (§4.5) will handle the conflict resolution logic to ensure convergence.

The challenge here is that Mnesia, by default, considers dead nodes not to be part of the cluster and carries on operations as if they

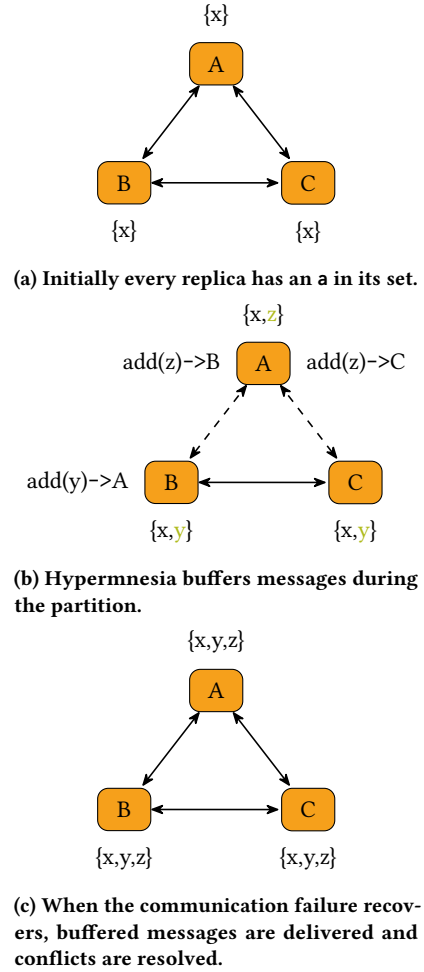


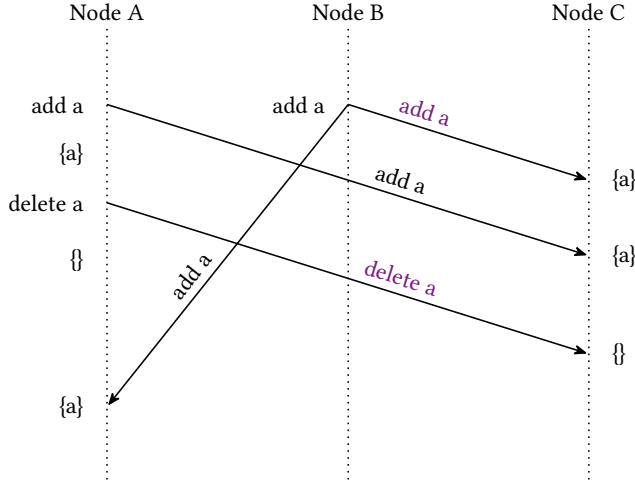
Figure 5: Hypermnesia buffers messages during communication failure and resolves conflicts afterwards.

did not exist. This might be a desirable behaviour in transactions, and developers can use the majority option to protect mission-critical data. But in an eventually consistent system, this is less ideal since we want our system to repair itself automatically. Therefore Hypermnesia considers dead nodes as part of the cluster and buffers operations for them.

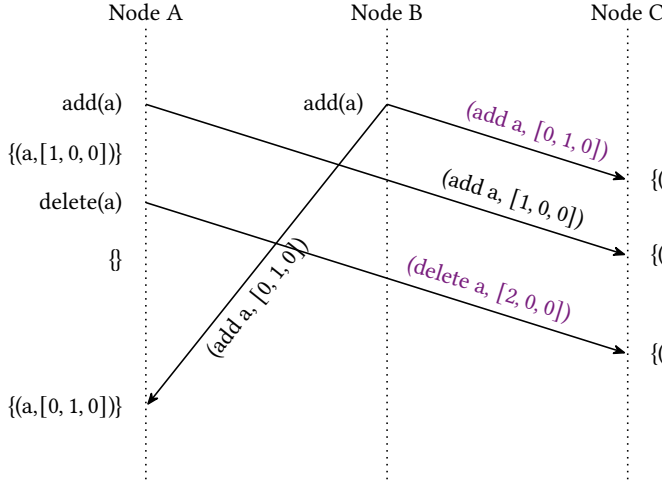
4.12 Transient failure

When a transient network failure happens, communication between nodes temporarily stops but is not long enough for the failure detector to act. Transactions will completely stall during this period. Although dirty operations can carry on, replicas might end up in different states due to out-of-order message delivery. For example, in fig. 6a, there might be a network failure between node B to A, resulting in B's add a being delayed. If messages are delivered as they arrive, then node A and node C will end up in an inconsistent state. This is because addition and deletion in a set do not commute, and the two purple operations (add and delete) are concurrent, and

they are applied in a different order on node A and node C, resulting in different final states.



(a) Mnesia replica states can diverge depending on the ordering of delivery of operations. The purple operations are causally concurrent.



(b) Hypermnnesia attaches vector timestamps to each message and stores them along with actual elements in the set for conflict resolution.

Figure 6: Mnesia and Hypermnnesia’s response to transient network failure.

In order to achieve convergence, and hence eventual consistency, concurrent operations need to commute. The exact semantics of whether addition or deletion wins depends on the actual application, and to achieve convergence, it is sufficient to define consistent semantics across replicas. Add-wins semantics is presented here but the remove-wins semantics is similar (??). To achieve add-wins semantics with a pure op-based Set CRDT, we require deletions to only remove elements that causally precede it, as captured by ?? 2 in alg. 2. In fig. 6b, the deletion with timestamp [2, 0, 0] removes

the element $(a, [1, 0, 0])$ which is causally lower, but not $(a, [0, 1, 0])$ which is causally concurrent.

4.13 Summary

This chapter covers the design of Hypermnnesia (??) in terms of its APIs to minimise the amount of code refactoring needed (§3.1). We also examined the suitability of different CRDTs and decided to use pure op-based CRDTs primarily because of its similarity to Mnesia’s communication pattern (§3.2). A prototype implementation is presented (§4.1), including a Tagged Causal Stable Broadcast layer by extending the basic causal broadcast with timestamp and causal stability information (§4.4). Practical optimisations such as the use of bag data structures are also considered (§4.5). We also briefly talked about two patterns of network faults and how the AW-set and buffering of operations can help us (§4.10).

5 EXPERIMENTAL EVALUATION

This chapter evaluates Hypermnnesia against the proposed research questions (crefsec:intro motivation). We start by introducing new tests to ensure the correctness of the system (§5.2), and then outline the benchmarking approach (§5.3), before performing measurements on a cluster of distributed physical machines in terms of time (§5.3) and space (§5.9). This is followed by an examination of the fault tolerance properties of Hypermnnesia (§5.10). We conclude the evaluation with a discussion on the API usability of Hypermnnesia (§5.11).

5.1 Setup

The experiments are performed on: (a) a single physical machine with multiple Erlang VMs connected via the loopback interface, mainly used for correctness testing and some simple experiments; (b) multiple physical machines inside a data centre connected via ethernet switches, used for more extensive benchmarking.

The single physical machine has an AMD 12-Core Processor and 16GiB DDR4 memory. It runs the Ubuntu 22.04.2 LTS Operating System. The cluster consists of 10 machines, each has an Intel(R) Xeon(R) CPU with 6 cores and 64GiB memory. They all run the Ubuntu 20.04.4 LTS Operating System. More details on the test-bed setup can be found in ??.

5.2 Correctness

The first question of whether eventual consistency is possible in Mnesia (item RQ1 in §1.1) is concerned with the design and correctness of the implementation, therefore correctness testing is performed. The Mnesia source code has a regression test suite with about 5000 tests, covering various aspects of the correctness such as transaction, dirty operations, storage engine, etc. They are run against the additional code introduced by Hypermnnesia. These mostly serve as tests for backwards compatibility, since there is currently no test that simulates network partitions, nor tests targeting the new eventually consistent API.

To cover the new EC API, the suite is extended by adding the following tests:

- (1) Unit tests on the behaviour of the AW-set and causal delivery, such as testing, among others, whether the addition of an element is reflected in a subsequent read.

- (2) Unit Tests on database features such as index reads.
- (3) Integration tests for simulating network partitions, including multiple cases of concurrent addition and deletion.

There are about 20 new tests written to cover the new API. Both unit tests and integration tests generally follow the pattern of the Erlang EUnit and Common Test framework [14, 16], respecting the existing test suite structure. The network partition is simulated with a custom Erlang distribution protocol [17], developed by RabbitMQ for their integration testing [45].

Hypermnnesia passes all of the tests mentioned above, and the output log of the tests is included in ??.

5.3 Benchmarking

5.4 Benchmark overview

The second question (item RQ2 in §1.1) is concerned with the efficiency of the eventual consistency API in Hypermnnesia. To answer this question, benchmarking is performed on the Hypermnnesia codebase. The original Mnesia repository has a built-in benchmarking suite for transactions. Briefly, this benchmark is a simulation of workloads where Mnesia tables are used to store session data of, for example, users logging into a website. There are four stages in this benchmark:

- population** initialises tables with randomly generated data.
- warmup** performs operations to bring data from memory to the cache.
- actual benchmarking** performs the actual benchmarking.
- cooldown** allows the system to clean up resources and ensures isolation between consecutive runs.

Typical operations in this benchmark include reading data of a particular session, creating a new session for a subscriber of the service, deleting the details of a particular session, etc.

To compare different access contexts, the benchmark is extended from transactions to dirty and EC operations. The extension is a substitution of transaction access context for dirty/EC access contexts. One thing to note is that, unlike transactions, ACID properties are not guaranteed for dirty and EC operations. Therefore it is common to have failures while reading data as messages might still be in transit. In such cases, we retry several times before considering the operation as failed and aborting it.

The general benchmarking strategy in the following sections is to measure the throughput and latency across three access contexts: dirty, transaction and EC, and make comparisons between them. Several metrics are considered in this benchmarking process: number of generators, number of nodes, table size, and workload types (§§5.5 to 5.8). We expect dirty operations to be the most performant in terms of throughput and latency since it has the lowest operational complexity, while transactions to be the slowest due to its synchronisation overhead. EC operations should stay between those two, and ideally as close to dirty operations as possible.

5.5 Number of generators

Figure 7 shows the comparison of throughput and latency against the number of generators per node for different access contexts. Generators are clients sending requests to the Mnesia cluster. These

requests are sent to the nearest node (with respect to the generator) for processing.

Generally speaking, throughput (fig. 7a) increases as the number of generators per node increases. For dirty operations, this saturates at around four generators per node, likely caused by I/O bandwidth limitation. Transactions and EC throughput continue to increase as more generators are added. This scalability property likely comes from the fact that Mnesia is, by default, a leaderless architecture (§2.4) where every node can process requests. Such design choices help Mnesia obtain more parallelism as we add more nodes into the cluster. Latency (fig. 7b), on the other hand, generally stays the same. Dirty operations do show a higher increasing rate than the other two, possibly because the overhead of more generators (and hence more messages to process) are more significant in a previously low-latency environment (the latency of dirty operations is on the order of 10 μ s).

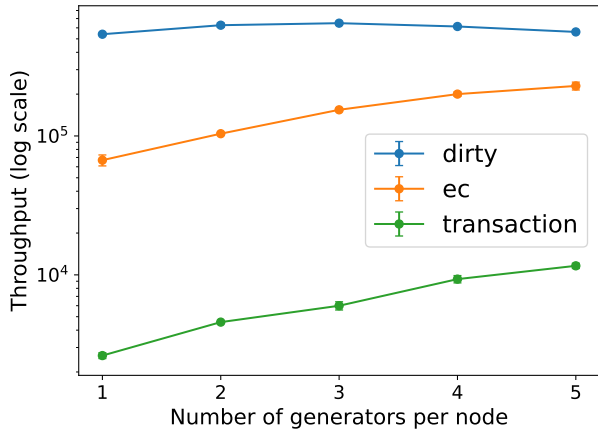
The increasing throughput and stable latency of EC operations demonstrate its scalability against the number of clients. Moreover, as the number of generators per node increases, the throughput of EC operations approaches about half the throughput of dirty operations, whereas initially it was only about one-tenth. This is a desirable property in a replicated system as it is scalable with respect to the number of clients (requests) it can handle.

5.6 Cluster size

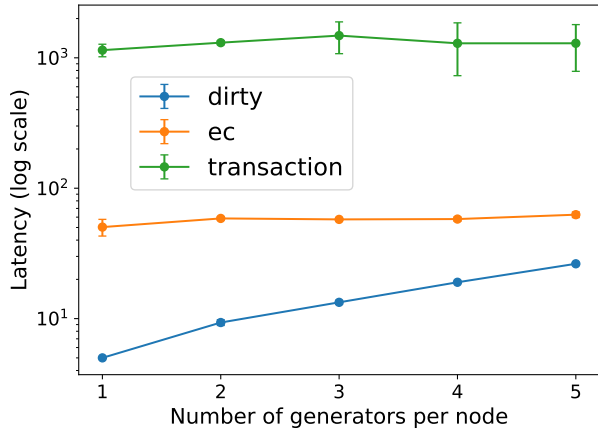
This section evaluates the throughput and latency against the number of nodes in the Mnesia cluster. Based on the original design of Mnesia (§2.2), the number of nodes used in a cluster generally does not exceed ten [27]. Figure 8a and fig. 8b show the throughput and latency change with respect to the number of nodes, with a fixed number (2) of generators. We observe that the throughput decreases and the latency increases as the number of nodes increases, suggesting that adding more nodes inevitably introduces overhead into the system, e.g. more messages to send, and more data to process. Keeping the number of generators the same means we keep the total amount of client requests the same but add duplicate work by adding more replicas, therefore the system experiences more overhead.

However, we could also increase the number of generators as we have more replicas, which is generally what happens in a real deployment as the number of nodes in a cluster of replicated machines increases. Figure 8c and fig. 8d show the throughput and latency as the number of generators increases linearly with the number of nodes. For all three operations, throughput increase is observed. Dirty operations' throughput saturates at about 7 nodes, when the message queue backlog starts to become the bottleneck. In terms of latency, there is an increase in all three. EC operations demonstrate around 13x higher throughput than transactions when there are nine nodes.

A higher replication factor often implies extra work for the system, thus resulting in lower throughput and higher latency. However, we could offset this with more clients and hence more parallelism. The overall effect is an increasing throughput as the number of nodes increases, albeit with an inevitable sacrifice in latency. For EC operations, this is approximately 60 μ s. This property is present in all three access contexts.



(a) Throughput



(b) Latency

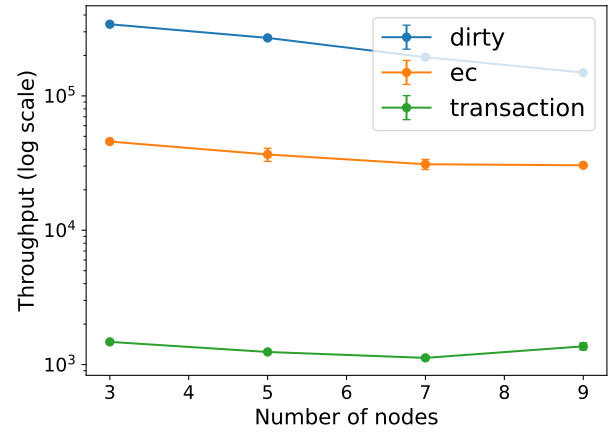
Figure 7: Throughput and latency against the number of generators. More generators generally bring higher throughput and latency.

5.7 Workload types

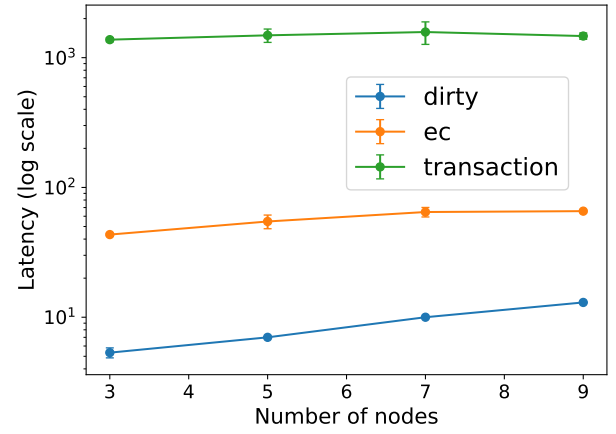
Figure 9 compares the throughput and latency of three access contexts against different workloads. Dirty operations are about 50x better than transactions, while EC operations lie between them (again), with about 10x higher throughput than transactions.

Note that if the read percentage is 100%, i.e. a read-only workload, the transaction gets close to or even surpasses the performance of EC operations. This is due to an optimisation done in the benchmark where it uses a synchronous dirty operation for reading data rather than a full two-phase commit, which removes the cost of locking, and multi-round communication time. It is even faster than EC operations as it does not need to go through the causal broadcast layer and the processing logic of an AW-set.

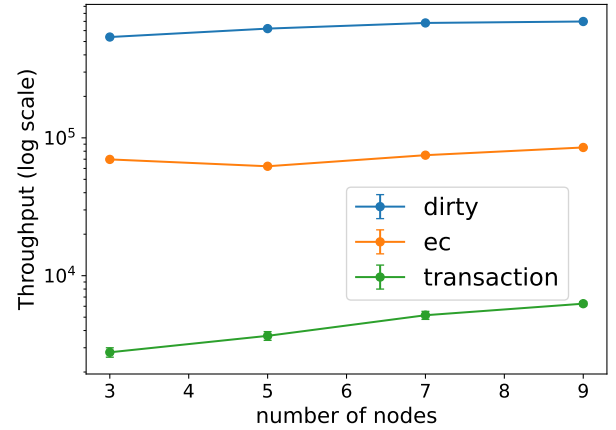
Write-intensive workloads are generally slower than read-intensive ones, since writes need to involve all replicas in the cluster, while reads can be done locally (if the data is present on the local node, which is the case in this benchmark). This trend is, again, true for



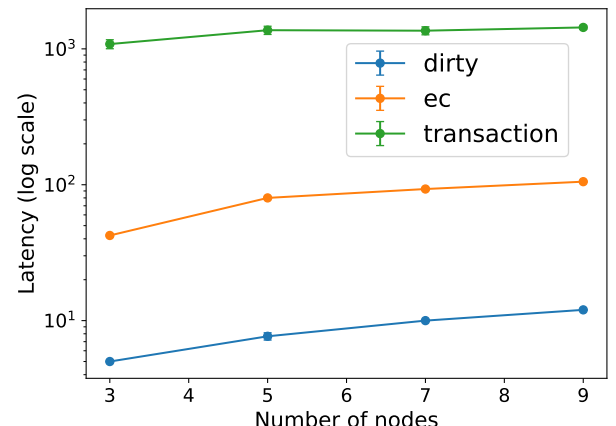
(a) Throughput with a fixed number of generators.

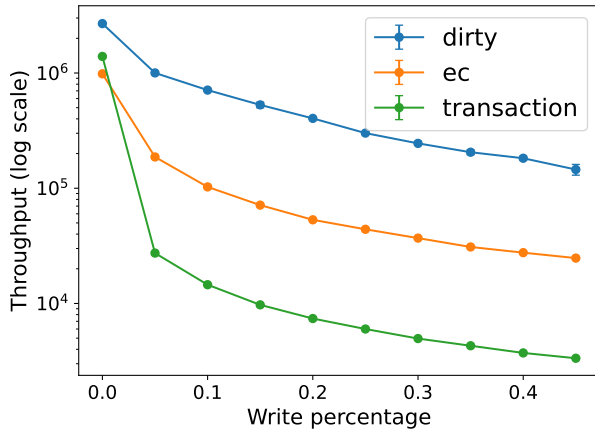


(b) Latency with a fixed number of generators.

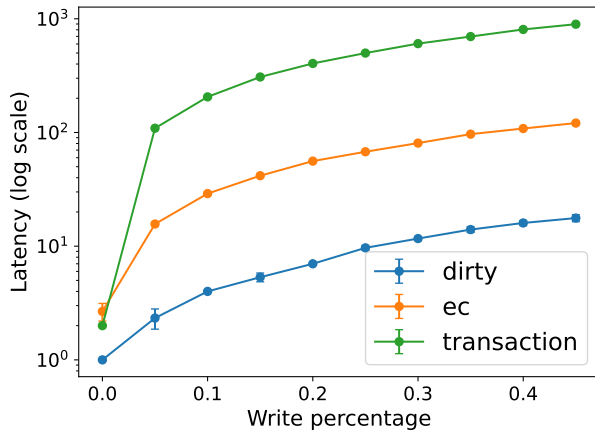


(c) Throughput with varying number of generators.





(a) Throughput



(b) Latency

Figure 9: Throughput and latency against different workload types. Read-heavy workloads are faster than write-heavy ones in all three cases.

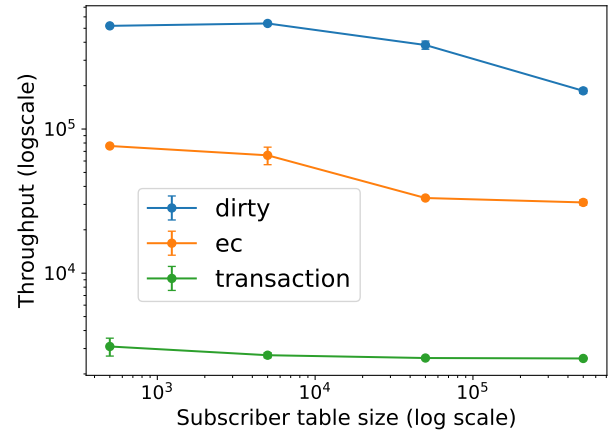
all three contexts. As long as the workload is not read-only, EC operations performs better than transactions.

5.8 Table size

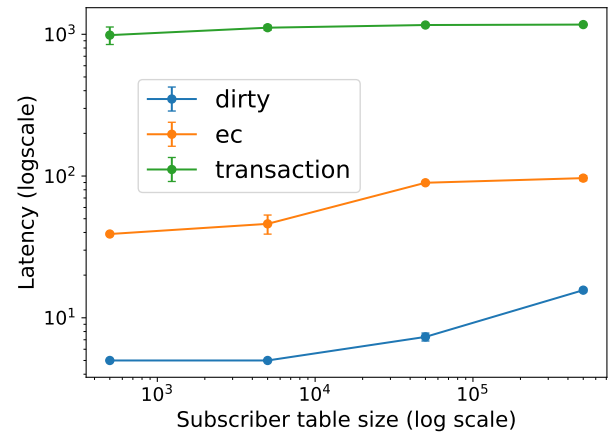
Figure 10 shows the change in throughput and latency as we vary the size of the subscriber table. There are five tables in total in this benchmark, and the subscriber table stores the data for users subscribing to a service, which is the most frequently changed and the largest one among all five tables. Therefore I choose to vary the subscriber table size during the initial table population.

A larger table generally makes it slower to read/write data, which affects dirty operations but has less impact on transactions and EC operations. Similar to the behaviour in §5.5, the low-latency dirty operations are more sensitive to even small changes in the latency of each operation, and exhibits a larger increase in latency.

Table size’s impact on the performance of Hypermnesia mainly comes from the extra time in accessing elements of a larger table.



(a) Throughput



(b) Latency

Figure 10: Throughput and latency against table size. Dirty operations are affected the most.

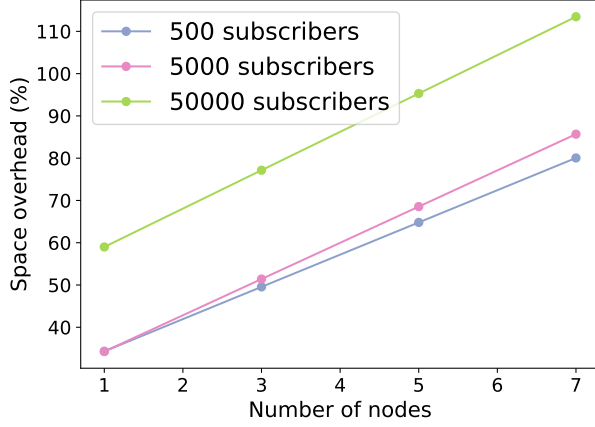
For EC operations, the extra overhead of periodic cleaning of timestamps could also play a role (§4.7). However, this overhead is still acceptable as EC operations still have about 25x higher throughput than transactions.

5.9 Space overhead

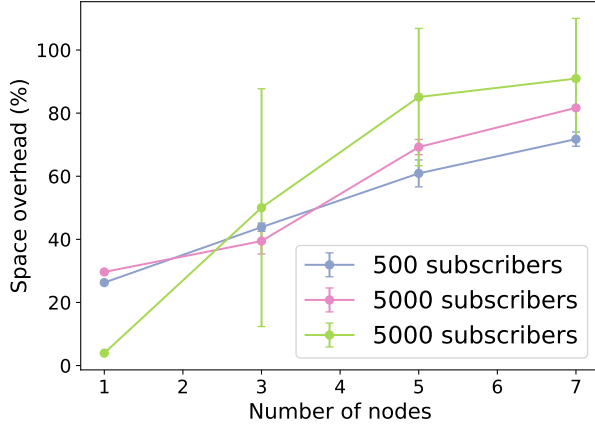
This section evaluates the space overhead of Hypermnesia. Generally speaking, CRDTs rely heavily on metadata to keep track of the history of operations and hence has a large overhead in its space usage [7]. In the case of implementing a pure AW-set, the primary metadata associated with each element is the vector clock timestamp, which grows linearly with the number of nodes in the cluster, although garbage collection is implemented to reduce its impact (§4.5).

Figure 11a shows the pure op-based set’s overhead compared to the default set in Mnesia. The overhead grows linearly with respect to the number of nodes as expected. Note that this is a static process, i.e. each time the populator will put the same number of elements into the set, therefore the lines are perfectly straight

with no variance. Figure 11b shows the overhead after running the benchmark. This is a dynamic process during which the causal stability optimisation is applied, but not in the previous figure (we discuss the reason below). Observe that the causal stability optimisation is able to reduce the overhead by up to 30%.



(a) Space overhead after populating the table.



(b) Space overhead after running the benchmark.

Figure 11: Space overhead. The number of subscribers represents the size of the table.

Despite space optimisations, the space overhead is still relatively large, especially when the number of nodes exceeds five. There are several reasons for this:

- The underlying implementation of the PO-Log is a set-like structure rather than a partially ordered log (§ 4.9). This complicates the implementation of causal stability optimisation as it is now harder to find all elements with a timestamp smaller than the stable one.
- The causal stability optimisation relies on nodes continuously receiving updates from other nodes to determine timestamp stability. This is not always possible since it is common for a node to only receive messages from others (maybe no client sends requests to it). This is why the optimisation is not applied in the population phase. Bauwens

and Gonzalez Boix [7] made a similar observation and proposed a solution based on an eager collection of metadata.

Memory management has been an important issue in CRDT research but has yet to attract much attention [7]. The current optimisation using causal stability is less effective than one would hope for. In a typical setup of three nodes, the extra space needed is around 30–40%. For this reason, Hypermnnesia is limited to relatively small clusters, but this is acceptable since Mnesia is designed for small clusters in the first place [27]. We discuss more on how to optimise the space overhead of the Set CRDT in § 7.2.

5.10 Fault tolerance

§ 4.10 discusses how Mnesia and Hypermnnesia respond to network partitions. In this section, we evaluate Hypermnnesia against these two situations to answer the research question item RQ3 in § 1.1.

Communication failure. Mnesia does not handle communication failure by default and asks application developers to resolve conflicts. Hypermnnesia buffers operations during the partition until it recovers (fig. 5c). It then sends the buffered message and resolves conflicts automatically.

Transient failure. As discussed in § 2.7, during a transient failure, transactions stall until the partition heals. We measure this effect by examining the throughput change in a simulated network partition, as shown in fig. 12. Figure 12a shows how the throughput changes through time, which is relatively stable when there is no partition. On the other hand, fig. 12b shows the throughput when there is a (simulated) network partition, lasting around 10 seconds. Dirty and EC operations remain about the same, but transaction throughput drops down to zero during the partition period. This is expected since the two-phase commit protocol requires a reply from *all* participating nodes.

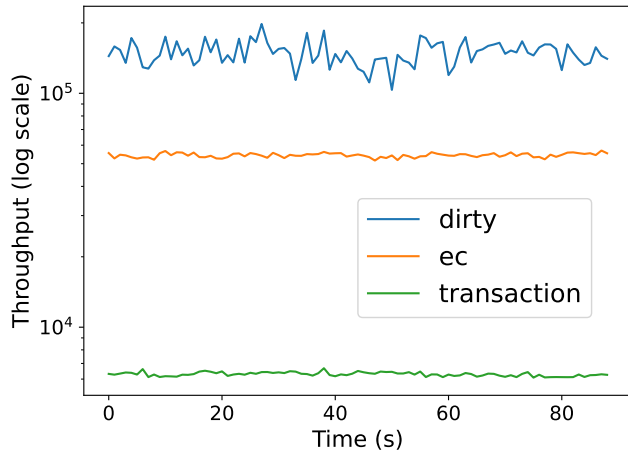
In summary, Hypermnnesia adds fault tolerance to Mnesia in the presence of network partitions: (1) automatic conflict resolution after a communication failure; (2) nodes remain available during transient failure, with guaranteed convergence when the partition heals.

5.11 APIs and refactoring

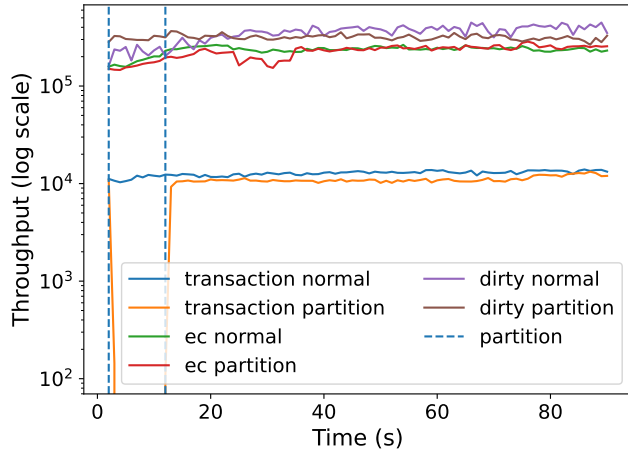
This section attempts to answer the research question item RQ4 concerning refactoring. First, a summary of the code changes is given for Hypermnnesia, followed by experimental modifications of real-world applications to demonstrate Hypermnnesia’s usability.

Listing 10 shows the code changes needed from transaction-/dirty operations (first two columns) to EC operations (last column). Note that refactoring is often just one line of change. Another change is the type option while creating a Mnesia table, as shown in listing 11. And that is all the change needed to use the new API.

The second step is to apply the refactoring described above in actual production codebases such as RabbitMQ [55] and ejabberd [43], and run the regression test suites against such changes. To be more cautious while refactoring, most changes are from asynchronous dirty operations to asynchronous EC operations. Hypermnnesia successfully passes the test suites of both applications with little effort in changing the source code.



(a) No partition.



(b) A partition of 10 seconds, dashed lines indicate the start and end of the partition.

Figure 12: Throughput changes against time. Network partition affects the transaction throughput but not dirty and EC throughputs.

It is worth noting that this is not a formal usability study of the new API. A more rigorous study requires much more extensive testing and understanding of the project codebases, which requires lots of engineering effort and is beyond the scope of this project. Nevertheless, we believe this evaluation is a first step towards making Hypermnesia more usable and production-ready.

5.12 Summary

In this chapter, we evaluated Hypermnesia against the research questions of this project (§1.1). The results show that Hypermnesia can produce correct results in spite of network delays (§5.2), and partitions (§5.10), thus demonstrating the possibility of eventual consistency in Mnesia and its role in automatic conflict resolution (items RQ1 and RQ3). Furthermore, benchmarking results show that EC operations can achieve approximately 10–20x better throughput

```
mnesia:activity(transaction,
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).

mnesia:activity(async_dirty,
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).
```

(a) transactions with activity/2. (b) dirty with activity/2.

```
mnesia:activity(async_ec,
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).
```

(c) EC operations with activity/2

```
mnesia:transaction(
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).

mnesia:async_dirty(
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).
```

(d) transactions with transaction/1 (e) dirty with async_dirty/1

```
mnesia:async_ec(
fun () ->
  mnesia:write(tab, tup),
  mnesia:read(tab, key)
end).
```

(f) EC operations with async_ec/1

Listing 10: Changing from transaction or dirty operations to EC operations.

```
ejabberd_mnesia:create(?MODULE, oauth_client, [{disc_copies,
  [node()]},
  {attributes, record_info(fields, oauth_client)}, {type, set}])).
```

(a) Original code creating a Mnesia table, using a set data structure.

```
ejabberd_mnesia:create(?MODULE, oauth_client, [{disc_copies,
  [node()]},
  {attributes, record_info(fields, oauth_client)}, {type, pawset}])).
```

(b) New code using the pure AW-set (pawset).

Listing 11: Adding a type declaration when creating a Mnesia table. Code excerpt modified from ejabberd [43].

and lower latency than transactions, and its performance can get close to dirty operations (recall that dirty operations are much faster than transactions in Mnesia from §2.6) when increasing the scale of the experiment (§5.3). This makes EC operations competitive for real-world applications (item RQ2). Finally, §5.11 evaluates the usability of the new API (item RQ4) and shows that Hypermnesia’s API enables minimum code refactoring for adoption in real-world projects.

6 RELATED WORK

The field of database research is a dynamic and broad domain that encompasses a variety of topics. This chapter focuses on the

architecture and design of eventually consistent databases (§6.1). Similarly, for CRDTs (§6.2), this chapter surveys them in the context of how they impact the design and implementation of Hypermnesia.

6.1 Databases and eventual consistency

6.1.1 Mnesia and eventual consistency. Mnesia’s lack of automatic conflict resolution has been problematic for developers using it [40, 60]. Unsplit [59], an external library for Mnesia, allows for user-defined merge logic but requires developers to define the custom logic. ForGETS [57] is WhatsApp’s Mnesia “drop-in” replacement database, aiming to provide Mnesia’s missing features such as auto reconciliation and auto reconnection. ForGETS uses the last-write-wins strategy for conflict resolution. However, ForGETS is a proprietary database influenced by WhatsApp’s operational experience, and the details of its effectiveness are unknown. Moreover, ForGETS is a standalone database built on top of ets, while Hypermnesia aims to extend Mnesia to support automatic resolution natively.

6.1.2 Key-value stores. Key-value stores like Redis³ serve as a caching layer between servers and databases, offering an Active-Active architecture that allows replicated database instances to distribute over different (possibly distant) locations. Redis is an advanced system with many more features such as expiring keys, and uses op-based CRDTs for conflict resolution between different instances [46]. While Mnesia has a similar use case, it is primitive and embedded into the Erlang ecosystem. Mnesia can provide a faster response time than Redis when used on a smaller scale, thanks to its shared address space with the application. This research aims to enhance Mnesia with stronger non-transactional consistency models found in a general purpose in-memory database (like Redis), reducing the constraints when developers opt for Mnesia as the choice of their Erlang applications.

6.1.3 Multi-version concurrency control. SwiftCloud [42] is a fault-tolerant geo-replicated transactional system. It provides causally consistent snapshots for each transaction with object versioning. As a geo-replicated system, SwiftCloud allows local operations to be performed on the client without going through the server but only when they operate on mergeable data types (e.g. CRDTs) during a transaction. Antidote database [49] goes one step further, integrating TCC with stronger consistency models and allowing developers to choose between them.

Multi-version Concurrency Control (MVCC) is a popular technique used in conjunction with CRDTs to provide eventual or sometimes even transactional consistency. Mnesia does not support object versioning by default, so adding MVCC requires extensive work. In this project, we focus on adding eventual consistency and leave the integration of MVCC and CRDTs as future work (§7.2).

6.1.4 Augmenting existing embedded databases. SQLite⁴ is one of the most widely used embedded SQL database engines in the world [28]. It does not have built-in support for replication, but there has been work that extend it for local-first software [52] with Conflict-free replicated relations (CRRs): CRDTs applied to relational databases [64]. They use a two-layer architecture: an application layer for handling client requests, and a CRR layer for

conflict resolution and anti-entropy protocols. However, their work still remains in progress and it is not clear how effective their approach is in terms of overheads. ElectricSQL⁵ is another attempt to augment SQLite with local-first behaviours. It uses RichCRDT [3], which are CRDTs extended with database guarantees, based on ideas from Antidote.

These techniques for enhancing SQLite are similar in spirit to Hypermnesia. However, SQLite is a rather general purpose single-node database receiving constant developments. In contrast, Mnesia is integrated with the Erlang ecosystem and designed as a distributed database, hence they have different use cases. Moreover, to the best of our knowledge, there is no performance evaluation on these extensions to SQLite yet.

6.2 CRDTs research

6.2.1 Systems using CRDTs. CRDTs are often used in collaboration software such as shared text editors [58] to provide local-first behaviours [31], i.e. users of the editors can keep typing into the document despite unstable networks. It has then found its use in areas such as synchronisation on the edge and opportunistic networks [25, 63]. These systems typically have constrained network or computing resources and CRDTs can be used to delay synchronisation and respond to the user first.

Researchers have been investigating new areas where CRDT techniques can be applied. Hypermnesia does not aim to apply CRDTs in a novel way but rather to experiment with their suitability in an embedded database like Mnesia.

6.2.2 Time and space improvements. CRDTs are constantly invented and improved to reduce their resource consumption. For example, the proposal of δ -CRDTs [1] is partly due to the communication overhead of propagating the entire state in normal state-based CRDTs. Pure op-based CRDTs [4] were also proposed to reduce the space overhead. Bauwens and Boix [6] further improve the reactivity and storage requirements of pure op-based CRDTs by exposing more information from the causal broadcast layer. van der Linde et al. [53] suggest a new way to adapt δ -CRDTs for more unstable networks. Furthermore, Bauwens and Gonzalez Boix [7] show a more eager way to remove metadata using acknowledgements from replicas.

One research question of this project is to understand the overhead of eventual consistency in Mnesia. The space and time optimisations proposed are useful for improving Hypermnesia’s efficiency, therefore they are investigated in detail, and some of their ideas are used in Hypermnesia’s implementation (§4.5).

7 CONCLUSION

In this project, we looked at Hypermnesia, an extension to the Mnesia DBMS incorporating eventual consistency. We conclude this report by summarising the accomplishments of Hypermnesia, highlighting key contributions (§7.1) and pointing out ways Hypermnesia can be improved in the future (§7.2).

³<https://redis.com>

⁴<https://sqlite.org/index.html>

⁵<https://electric-sql.com>

7.1 Accomplishments

Qualitatively speaking, Hypermnesia’s API is designed to minimise the amount of code refactoring (??) and fits well into the existing Mnesia access contexts, making it easier to be adopted in existing open source codebases (§5.11) such as RabbitMQ and ejabberd [43, 55]. Moreover, the new eventual consistency API passes the extended Mnesia regression test suite (consisting of ≈ 5000 unit tests), including (≈ 20) additional tests that cover network partition, AW-set and RW-set behaviours and causal broadcast (§5.2). It also allows the database to continue operating while there is a partition and automatically resolves potential conflicts after the partition recovers (§5.10).

Quantitatively speaking, Hypermnesia accomplishes the above functionalities while providing about 10–20 times higher throughput and lower latency than Mnesia’s transactions (§5.3). Thanks to Mnesia’s performant dirty operations, Hypermnesia can be designed to guarantee eventual consistency without sacrificing much performance, taking advantage of the performant architecture of dirty operations.

7.2 Future work

Finally, I list some possible extensions to Hypermnesia:

- At the moment EC operations do not interact well with Mnesia’s transactions and dirty operations, i.e. they cannot be used on the same table. As mentioned in §6.1.3, it is possible to add support for transactions that execute and update queries on a consistent snapshot and then resolve conflicts between different snapshots with CRDTs, and there are many protocols developed for this purpose [42, 49]. Now that Mnesia has built-in support for CRDTs and eventual consistency, it could be further enhanced to support lightweight but highly available transactions.
- Hypermnesia currently works for in-memory tables only. Although in-memory caching tends to be the way Mnesia is used [37], Hypermnesia could be extended to support disk tables in the future, or even custom backends. Different data structures might open the door for better space optimisation which is currently less feasible with ets and dets.
- Hypermnesia chooses to use pure op-based CRDTs for its simplicity and similarity to Mnesia’s architecture. δ -CRDTs is another popular choice among eventually consistent databases [32], which might be worthwhile experimenting with.

ACKNOWLEDGMENTS

This work was supported by the [...] Research Fund of [...] (Number [...]). Additional funding was provided by [...] and [...]. We also thank [...] for contributing [...].

REFERENCES

- [1] Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero. 2018. Delta State Replicated Data Types. *J. Parallel and Distrib. Comput.* 111 (Jan. 2018), 162–173. <https://doi.org/10.1016/j.jpdc.2017.08.003> arXiv:1603.01529 [cs]
- [2] Jesper L. Andersen. 2014. Mnesia and CAP.
- [3] Valter Balegas. 2022. Introducing Rich-CRDTs - ElectricSQL.
- [4] Carlos Baquero, Paulo Sérgio Almeida, and Ali Shoker. 2014. Making Operation-Based CRDTs Operation-Based. In *Proceedings of the First Workshop on Principles and Practice of Eventual Consistency (PaPEC ’14)*. Association for Computing Machinery, New York, NY, USA, 1–2. <https://doi.org/10.1145/2596631.2596632>
- [5] Carlos Baquero, Paulo Sérgio Almeida, and Ali Shoker. 2017. Pure Operation-Based Replicated Data Types. <https://doi.org/10.48550/arXiv.1710.04469> arXiv:arXiv:1710.04469
- [6] Jim Bauwens and Elisa Gonzalez Boix. 2021. Improving the Reactivity of Pure Operation-Based CRDTs. In *Proceedings of the 8th Workshop on Principles and Practice of Consistency for Distributed Data*. ACM, Online United Kingdom, 1–6. <https://doi.org/10.1145/3447865.3457968>
- [7] Jim Bauwens and Elisa Gonzalez Boix. 2019. Memory Efficient CRDTs in Dynamic Environments. In *Proceedings of the 11th ACM SIGPLAN International Workshop on Virtual Machines and Intermediate Languages (VMIL 2019)*. Association for Computing Machinery, New York, NY, USA, 48–57. <https://doi.org/10.1145/3358504.3361231>
- [8] David Bernbach and Jörn Kuhlenskamp. 2013. Consistency in Distributed Storage Systems. In *Networked Systems (Lecture Notes in Computer Science)*, Vincent Gramoli and Rachid Guerraoui (Eds.). Springer, Berlin, Heidelberg, 175–189. https://doi.org/10.1007/978-3-642-40148-0_13
- [9] Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman. 1987. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley Pub. Co, Reading, Mass.
- [10] Kenneth Birman, André Schiper, and Pat Stephenson. 1991. Lightweight Causal and Atomic Group Multicast. *ACM Transactions on Computer Systems* 9, 3 (Aug. 1991), 272–314. <https://doi.org/10.1145/128738.128742>
- [11] Francesco Cesarini. 2019. Companies Who Use Erlang.
- [12] Alan Demers, Dan Greene, Carl Hauser, Wes Irish, John Larson, Scott Shenker, Howard Sturgis, Dan Swinehart, and Doug Terry. 1987. Epidemic Algorithms for Replicated Database Maintenance. In *Proceedings of the Sixth Annual ACM Symposium on Principles of Distributed Computing (PODC ’87)*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/41840.41841>
- [13] C. A. Ellis and S. J. Gibbs. 1989. Concurrency Control in Groupware Systems. *ACM SIGMOD Record* 18, 2 (June 1989), 399–407. <https://doi.org/10.1145/66926.66963>
- [14] Ericsson AB. 2023. Erlang – Common Test Basics. https://www.erlang.org/doc/apps/common_test/basics_chapter.html
- [15] Ericsson AB. 2023. Erlang – Erlang Reference Manual. https://www.erlang.org/doc/reference_manual/users_guide.html
- [16] Ericsson AB. 2023. Erlang – EUnit - a Lightweight Unit Testing Framework for Erlang. <https://www.erlang.org/doc/apps/eunit/chapter.html>
- [17] Ericsson AB. 2023. Erlang – How to Implement an Alternative Carrier for the Erlang Distribution. https://www.erlang.org/doc/apps/erts/alt_dist.html
- [18] Ericsson AB. 2023. Erlang – Mnesia. <https://www.erlang.org/doc/man/mnesia.html>
- [19] Ericsson AB. 2023. Erlang – Mnesia User’s Guide. https://www.erlang.org/doc/apps/mnesia/users_guide.html
- [20] Ericsson AB. 2023. Erlang – OTP Design Principles. https://www.erlang.org/doc/design_principles/users_guide.html
- [21] Ericsson AB. 2023. Erlang – STDLIB User’s Guide. https://www.erlang.org/doc/apps/stdlib/users_guide.html
- [22] Erlang Solutions. 2023. *MongooseIM Platform*. Erlang Solutions.
- [23] Benjamin Farinier, Thomas Gazagnaire, and Anil Madhavapeddy. 2015. Mergeable Persistent Data Structures. In *Vingt-Sixièmes Journées Francophones Des Langues Applicatifs (JFLA 2015)*.
- [24] Seth Gilbert and Nancy Lynch. 2012. Perspectives on the CAP Theorem. *Computer* 45, 2 (Feb. 2012), 30–36. <https://doi.org/10.1109/MC.2011.389>
- [25] Frédéric Guédec, Yves Mahéo, and Camille Noûs. 2023. Supporting Conflict-Free Replicated Data Types in Opportunistic Networks. *Peer-to-Peer Networking and Applications* 16, 1 (Jan. 2023), 395–419. <https://doi.org/10.1007/s12083-022-01404-6>
- [26] Theo Haerder and Andreas Reuter. 1983. Principles of Transaction-Oriented Database Recovery. *Comput. Surveys* 15, 4 (Dec. 1983), 287–317. <https://doi.org/10.1145/289.291>
- [27] Fred Hébort. 2013. *Learn You Some Erlang for Great Good!: A Beginner’s Guide*. No Starch Press, San Francisco.
- [28] Richard Hipp. 2019. Most Widely Deployed SQL Database Engine. <https://www.sqlite.org/mostdeployed.html>
- [29] Martin Kleppmann. 2017. *Designing Data-intensive Applications: The Big Ideas Behind Reliable, Scalable, and Maintainable Systems*. O’Reilly Media.
- [30] Martin Kleppmann and Tim Harris. 2022. *Distributed Systems*. Cambridge CST Part IB Lecture Notes.
- [31] Martin Kleppmann, Adam Wiggins, Peter van Hardenberg, and Mark McGranaghan. 2019. Local-First Software: You Own Your Data, in Spite of the Cloud. In *Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward! 2019)*. Association for Computing Machinery, New York, NY, USA, 154–178. <https://doi.org/10.1145/3359591.3359737>
- [32] Rusty Klopheus. 2010. Riak Core: Building Distributed Applications without Shared State. In *ACM SIGPLAN Commercial Users of Functional Programming*

- (CUPP '10). Association for Computing Machinery, New York, NY, USA, 1. <https://doi.org/10.1145/1900160.1900176>
- [33] Leslie Lamport. 1978. Time, Clocks, and the Ordering of Events in a Distributed System. *Commun. ACM* 21, 7 (July 1978), 558–565. <https://doi.org/10.1145/359545.359563>
 - [34] Eran Levy. 2022. RabbitMQ vs Kafka: Use Cases, Performance & Architecture.
 - [35] Prince Mahajan, Lorenzo Alvisi, and Mike Dahlin. 2012. *Consistency, Availability, and Convergence*. Technical Report.
 - [36] Robert C Martin. 2000. *Design Principles and Design Patterns*.
 - [37] Håkan Mattsson, Hans Nilsson, and Claes Wikström. 1998. Mnesia — A Distributed Robust DBMS for Telecommunications Applications. In *Practical Aspects of Declarative Languages (Lecture Notes in Computer Science)*, Gopal Gupta (Ed.). Springer, Berlin, Heidelberg, 152–163. https://doi.org/10.1007/3-540-49201-1_11
 - [38] HAKAN MATTSSON. 1999. Mnesia Internals.
 - [39] HAKAN MATTSSON. 2009. Mnesia Implementation Documentation.
 - [40] Paul Mineiro. 2008. Dukes of Erl: Network Partition ... Oops.
 - [41] Nuno Preguiça. 2018. Conflict-Free Replicated Data Types: An Overview. arXiv:arXiv:1806.10254
 - [42] Nuno Preguiça, Marek Zawirski, Annette Bieniusa, Sergio Duarte, Valter Bolegas, Carlos Baquero, and Marc Shapiro. 2014. SwiftCloud: Fault-Tolerant Geo-Replication Integrated All the Way to the Client Machine. In *2014 IEEE 33rd International Symposium on Reliable Distributed Systems Workshops*. IEEE, Nara, Japan, 30–33. <https://doi.org/10.1109/SRDSW.2014.33>
 - [43] Processone. 2023. Processone/Ejabberd. ProcessOne.
 - [44] Mark Raasveldt and Hannes Mühleisen. 2019. DuckDB: An Embeddable Analytical Database. In *Proceedings of the 2019 International Conference on Management of Data (SIGMOD '19)*. Association for Computing Machinery, New York, NY, USA, 1981–1984. <https://doi.org/10.1145/3299869.3320212>
 - [45] RabbitMQ. 2022. Inet_tcp_proxy. RabbitMQ.
 - [46] Redis. 2022. Diving into Conflict-Free Replicated Data Types (CRDTs). <https://redis.com/blog/diving-into-crds/>.
 - [47] Jerome Saltzer and M. Frans Kaashoek. 2009. *Principles of Computer System Design: An Introduction*. Morgan Kaufmann.
 - [48] Frank B. Schmuck. 1988. *The Use of Efficient Broadcast Protocols in Asynchronous Distributed Systems*. Technical Report. Cornell University.
 - [49] Marc Shapiro, Annette Bieniusa, Nuno Preguiça, Valter Bolegas, and Christopher Meiklejohn. 2018. Just-Right Consistency: Reconciling Availability and Safety. <https://doi.org/10.48550/arXiv.1801.06340> arXiv:1801.06340 [cs]
 - [50] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. 2011. Conflict-Free Replicated Data Types. (2011).
 - [51] Dharma Shukla. 2018. Azure Cosmos DB: Pushing the Frontier of Globally Distributed Databases.
 - [52] Iver Toft Tomter and Weihai Yu. 2021. Augmenting SQLite for Local-First Software. In *New Trends in Database and Information Systems*, Ladjel Bellatreche, Marlon Dumas, Panagiotis Karras, Raimundas Matulevičius, Ahmed Awad, Matthias Weidlich, Mirjana Ivanović, and Olaf Hartig (Eds.). Vol. 1450. Springer International Publishing, Cham, 247–257. https://doi.org/10.1007/978-3-030-85082-1_22
 - [53] Albert van der Linde, João Leitão, and Nuno Preguiça. 2016. Δ -CRDTs: Making δ -CRDTs Delta-Based. In *Proceedings of the 2nd Workshop on the Principles and Practice of Consistency for Distributed Data*. ACM, London United Kingdom, 1–4. <https://doi.org/10.1145/2911151.2911163>
 - [54] Paolo Viotti and Marko Vukolić. 2016. Consistency in Non-Transactional Distributed Storage Systems. *Comput. Surveys* 49, 1 (June 2016), 19:1–19:34. <https://doi.org/10.1145/2926965>
 - [55] VMware. 2023. RabbitMQ Server. RabbitMQ.
 - [56] Werner Vogels. 2008. Eventually Consistent: Building Reliable Distributed Systems at a Worldwide Scale Demands Trade-Offs? Between Consistency and Availability. *Queue* 6, 6 (Oct. 2008), 14–19. <https://doi.org/10.1145/1466443.1466448>
 - [57] Mikhail Vorontsov. 2018. Mikhail Vorontsov - ForgeTS: A Globally Distributed Database - Code Beam STO.
 - [58] Stephane Weiss, Pascal Urso, and Pascal Molli. 2010. Logoot-Undo: Distributed Collaborative Editing System on P2P Networks. *IEEE Transactions on Parallel and Distributed Systems* 21, 8 (Aug. 2010), 1162–1174. <https://doi.org/10.1109/TPDS.2009.173>
 - [59] Ulf Wiger. 2023. Writing an Unsplit Method.
 - [60] Ulf Wiger. Thu Feb 4 22:39:02 CET 2010. [Erlang-Questions] Unsplit - Resolving Mnesia Inconsistencies.
 - [61] Wikipedia contributors. 2020. LYME (Software Bundle). *Wikipedia* (Dec. 2020).
 - [62] Wikipedia contributors. 2023. Eventual Consistency. *Wikipedia* (April 2023).
 - [63] Weihai Yu and Claudia-Lavinia Ignat. 2020. Conflict-Free Replicated Relations for Multi-Synchronous Database Management at Edge. In *2020 IEEE International Conference on Smart Data Services (SMDS)*. 113–121. <https://doi.org/10.1109/SMDS49396.2020.00021>
 - [64] Weihai Yu and Sigbjørn Rostad. 2020. A Low-Cost Set CRDT Based on Causal Lengths. In *Proceedings of the 7th Workshop on Principles and Practice of Consistency for Distributed Data (PaPoC '20)*. Association for Computing Machinery, New York, NY, USA, 1–6. <https://doi.org/10.1145/3380787.3393678>