

MASARYK UNIVERSITY
FACULTY OF INFORMATICS



Real-time collaboration in Komodo

MASTER'S THESIS

Bc. Matúš Makový

Brno, Spring 2015

Declaration

Hereby I declare, that this paper is my original authorial work, which I have worked out by my own. All sources, references and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Bc. Matúš Makový

Advisor: RNDr. Filip Nguyen

Acknowledgement

Abstract

Keywords

Contents

1	Introduction	1
2	Real-time collaboration	3
2.1	<i>Basic overview</i>	3
2.2	<i>Operational Transformation</i>	6
2.3	<i>Differential Synchronization</i>	15
2.4	<i>Commutative Replicated Data Types</i>	16
3	Comparison of Techniques	17
3.1	<i>Comparison Criteria</i>	17
3.2	<i>Comparison</i>	17
4	Komodo	18
4.1	<i>Overview</i>	18
4.2	<i>Komodo Requirements for Real-time Collaboration</i>	18
4.3	<i>Best Technique for Komodo</i>	18
4.4	<i>Java Implementation</i>	18
5	Conclusion	19

1 Introduction

Many software solutions enable people to create new things in a better and faster way. In most cases the resulting product should be so complex that one person is not enough for the successful and fast creation. Creators try to collaborate to achieve a common goal. Among other opportunities and possibilities are collaboration and sharing the greatest benefits of the Internet.

We can identify 2 types of collaboration over the Internet, non real-time collaboration and real-time collaboration. At the beginning, as the Internet didn't have such capacity, people tried to use it just for sharing their drafts of work and sending them to each other, this type of collaboration is called non real-time. In non real-time collaboration users work on separate copies of a project and then need to merge their changes into one final project. In other words, they had to find differences between their drafts and reflect them to each other's version. It doesn't offer such flexibility as real-time collaboration and also it had many limitations, for example 2 people could not edit the same file in a project without having to resolve conflicts manually when they tried to merge their work with collaborator's version

Examples of non real-time collaboration could be Revision control (Git, SVN). With the development of the Internet came a reasonable solution called real-time collaboration. Using this principle, author can see what his collaborator is doing in real-time and the manual synchronization or manual conflict resolution is not necessary. Information technologies take care of this synchronization and conflict resolution for the users.

In the theoretical part of this thesis, author deals with principles of real-time collaboration and describes some of the techniques used to implement real-time collaboration in software over the Internet. There is also a comparison of these techniques from various aspects.

The practical part of this thesis deals with JBoss Data Virtualization and Komodo software as a new version of Teiid Designer developed by Red Hat that should use real-time collaboration in its upcoming release. Author recommends best technique for this authoring software regarding the requirements of data models and finds a suitable implementation in Java programming language.

2 Real-time collaboration

This chapter covers the basic overview of collaboration in general, description of real-time collaboration and description of difficulties with its implementation over the Internet. Last sections of this chapter describe three techniques used for implementation over the Internet and its properties. Techniques described in this chapter are: Operational Transformation, Differential Synchronization and Commutative replicated data types.

2.1 Basic overview

Collaboration, as defined by English dictionary, is an act of working with another or others on a joint project. Authoring Systems in IT that support collaboration are called Groupware.

Groupware systems are computer-based systems that support two or more users engaged in a common task, and that provide an interface to a shared environment. These systems frequently require fine-granularity sharing of data and fast response times.[1]

There are many techniques used in groupware systems that are not suitable for real-time collaboration, for example, Locking or Single Active Participant technique. Fundamental principle of Locking and Single Active Participant is to lock data when it is being modified by someone. This principle is adopted for example by Microsoft, when users try to edit shared document.

When using real-time collaboration, authoring software creates an illusion that users are working on one common copy of a document online. There is no requirement to commit changes to some kind of shared repository and no need for a user to resolve conflicts. Changes should be reflected and saved immediately. Examples of real-time collaborative editors that we recognize today are Google Docs, Ether-

pad and already terminated Google Wave.

Software engineer has different options for implementation of real-time collaboration in software solution.

Requirements for a good technique are:

- speed
- latency toleration
- low data transfer
- consistency maintenance
- good conflict resolution

The speed requirement means that changes made on one side of the collaboration process need to be reflected to the other sides as soon as possible and vice versa. Changes should be sent to other collaborators as soon as they are done. If a collaboration technique is fast enough it is much easier to satisfy other requirements on this technique. Fast enough technique is able to maintain good consistency. The system's response time is the time necessary for the actions of one user to be reflected by their own interface; the notification time is the time necessary for one user's actions to be propagated to the remaining users' interfaces.[1] Response time and notification time should be as short as possible.

One of the problems of real-time collaborative editing could be the latency of network. Implemented technique should be able to tolerate the latency of internet connections, because collaborators could be in very distinct parts of the world. It should be able to reconstruct the right order of operations because data sent over the Internet don't necessarily come in right order and the order of the operations is very important to maintain consistency.

Low data transfer requirement is also very important. Different sides of the collaboration should send as few data as possible. Transferred data should only describe change that has been done on one side of collaboration, it is not necessary to transfer the whole project. The less data is needed to transfer the faster the whole protocol can be.

Consistency maintenance is necessary for the success, it has to be ensured, that users on both sides are looking at the same version of a document regardless of number and complexity of operations done on both sides of the collaboration. Lack of consistency could cause other problems and chaos in the document versions. According to [7] there are three problems encountered when trying to achieve consistency maintenance and they correspond to the properties of CCI consistency model proposed in [6]

1. **Casuality Preservation** - operation O_1 causally precedes operation O_2 if O_1 occurred locally before O_2 . The problem is to execute operations in right order on all sites.
2. **User Intension Preservation** - technique must preserve user's intension in the context of state, in which the operation was executed. This problem is in strong relation with conflict resolution requirement and occurs when it is not possible to determine if O_x causally precedes O_y or O_y causally precedes O_x .
3. **Convergence** - when same operations have been applied on every site, the documents are identical.

Because the real-time collaboration is asynchronous there arises a problem of concurrency. This means, that changes can happen at the same time and in the same sections of project. Good conflict resolution requirement is present because of concurrency. Implementation techniques have to be able to identify and resolve a conflict when users are editing the same part of the project.

The concurrency control algorithms are separated in two classes: pessimistic and optimistic.

Pessimistic algorithms require communication with other sites or with a central coordinator before making a change to data.[?] One example of pessimistic algorithm can be already mentioned Locking or Single Active Participant technique.

Optimistic concurrency control, on the other hand, requires no communication before applying changes locally. The party making a change applies it immediately, then informs the other parties of the action. If more than one participant makes a change at the same time, a conflict resolution algorithm creates compensating changes to move everyone to the same final state.[4]

This implies that optimistic algorithms are better solution for Internet, because the latency can not be guaranteed.

2.2 Operational Transformation

Operational transformation (OT) is optimistic concurrency control algorithm used for real-time collaborative editing over the Internet.

It was first introduced in paper Concurrency Control in Groupware Systems[1] in 1989, together with The Distributed Optional Transformation (dOPT) Algorithm. As the paper describes it, algorithm has a number of properties which make it suitable for groupware. First, operations are performed immediately on their originating site, thus responsiveness is good. Secondly, locks are not necessary so all data remains accessible to group members. Finally, the algorithm is fully distributed, and resilient to site failure.

This first algorithm was able to process only plain text. Later some problems with correctness were discovered and resolved in following works. Over the years many other algorithms implementing

OT have been published. Algorithms used today are able to process also XML and other formats.

OT was used also in Jupiter Collaboration System in 1995 [4]. Jupiter's two-way algorithm is derived from the dOPT algorithm used by Grove in [1]. This algorithm is one of the improvements of dOPT.

Today is Google one of the most common implementators of this protocol. The main application was in Google Wave project, which is now terminated, although Operational Transformation found application also in Google Docs project. According to [5] Google is one of the biggest innovators of this approach and its biggest contribution is the idea of operation composition, which is described later in this section.

The approaches differ also in the architecture, dOPT doesn't involve server in its design, but the recent algorithms use a central server that maintains its version of the document, coordinates the communication and broadcasts operations.

This section will explain basic principles of Operational transformation using Google version of the control algorithm, introduce some improvements added by Google and then point out the differences between the presented version and other versions like dOTP, Jupiter and TIPS.

The basic idea is that data are replicated on every machine and the only information that is sent over the Internet are the operations. Data replication ensures good responsiveness in high latency environments. Wave's addition to OT are also annotations. An annotation is some meta-data associated with an item range, i.e., a start position and an end position. This is particularly useful for describing text formatting and spelling suggestions, as it does not unnecessarily complicate the underlying structured document format.[8]

The main building block of this approach, as the name suggests, is

an operation. Here are operations defined for use in Google Wave[8]:

- `retain()` - move cursor
- `insertCharacters()` - inserts texts
- `insertElementStart()` - inserts starting tag
- `insertElementEnd()` - inserts end tag
- `deleteCharacters()` - deletes text
- `deleteElementStart()` - deletes starting tag
- `deleteElementEnd()` - deletes ending tag
- `replaceAttributes()` - replaces attributes in tag
- `updateAttributes()` - updates attributes in tag
- `annotationBoundary()` - describes changes in annotations

The operation is executed locally and then sent over the network to the server or to other peers, to execute the operation on their version of the document. Every object that can be changed in the document has its index for identification. Google Wave uses XML as format for storing information so as an object is considered a character in text and also a starting and terminating tag in XML structure (characters of tags are not considered elements).

The second major part of this technique is a transformation function, that transforms operations. Operations which are sent over the network and received by server can not be executed directly on the server's version of the document, because the index of the desired element could be different in the local copy and in the server copy of the document and this can cause inconsistency and violation of User Intension Preservation.

See Figure 2.1 with example for better explanation. For the purpose of this example, we will assume that we have one client that sends operations done on his local copy and a server that receives operations from all peers and maintains server version of the document.

Figure 2.1: Example Operational Transformation



Both of them start with text 'TAX' in their document. Pointer starts at index 0 and operation `retain()` shifts the pointer by number of indexes given as parameter. Client wants to add character 'E' so he executes operations:

```
retain(3);  
insertCharacters('E');
```

In the meantime server received operation from other peer, that has inserted character 'S', so server executes operations:

```
retain(3);  
insertCharacters('S');
```

Server and client exchange the operations. Different order and difference in indexes on both sides of the collaboration result in inconsistent document state. String in first document is 'TAXSE' and string in sec-

ond document is 'TAXES'. In order to preserve user's intension and to reach correct result we have to transform this operations. The correct result after inserting letters 'E' and 'S' should be 'TAXES'. Server has correct version of the document, but client doesn't so we have to transform the operation sent by server to look like this:

```
retain(4);
insertCharacters('S');
```

This transformation was helpful, but it is not enough. Now the client operation retains 3 characters, inserts new character and leaves cursor in front of last letter and the server operation retains 4 characters inserts character and leaves cursor behind the whole word. Client and server should end up in the same state. So it is necessary to transform also the client operation:

```
retain(3);
insertCharacters('E');
retain(1);
```

After this change the result of applying operations on both sides will be correct.

More gernerally, if we have two operations O_1 and O_2 we need a transformation function *trasform*, such that:

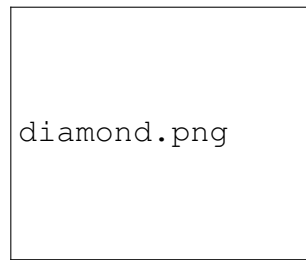
$$transform(O_1, O_2) = (O'_1, O'_2) \quad \text{where} \quad O'_2 \circ O_1 \equiv O'_1 \circ O_2$$

In other words, the function takes two operations and transforms them against each other and the output of the function are two transformed operations, such that, if we have same strings on both sides a we apply O_1 and O'_2 on one side, and O_2 and O'_1 on the other side, the resulting string is the same.

According to [3]. If an OT system is to support particular functionality, then it must be able to support certain transformation properties. For group editing and consistency maintenance, the system

must support a transformation function known as Inclusion Transformation. For group undo, where the effect of a previously executed operation is un-done at all sites, and all operations executed after it are all re-transformed, the system must support another transformation function known as Exclusion Transformation. Inclusion Transformation transforms operation O_A against another operation O_B in such a way that the impact of O_B is effectively included. Exclusion Transformation transforms operation O_A against another operation O_B in such a way that the impact of O_B is effectively excluded.[3] Problem described in the example could be also visualized as a "diamond problem". See Figure 2.2.

Figure 2.2: Diamond problem - 1



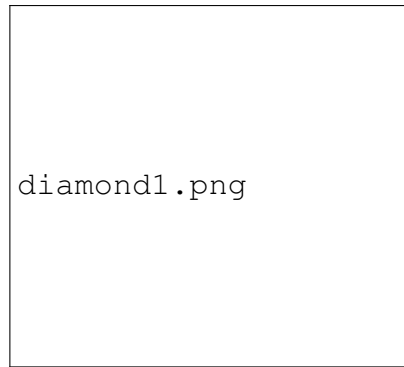
Operations done by peer take the document to the left side and operations done on server take the document to the right side. Each node in the diamond diagram represents the number of operations done by client and server in particular state.

As for the example, the correct scenario would be that peer and server exchange their operations, each one computes its transformation and applies the transformed operation on their versions of documents.

The problem presented by this example is actually the only one that is Operational Transformation capable of solving. Complicated situation occurs when sides of the collaboration diverge by more

then one step.

Figure 2.3: Diamond problem - 2



This example is used in [5] On the server side, transformation for operations a, c has to be computed, and server applies operation a' on its version of the document. The remaining operation c' has to be preserved for the next transformation. In the next step server must use b and c' as an input and computes b' that can be applied on its version of the document. Server's document is in the desired state.

On the client side, c has to be transformed against two operations to get to desired state. For two operations it is not a big problem, but as the number of operations could very quickly increase, because Operational Transformation is an optimistic technique, so it could potentially take long time to process. One of the advantages of Google Wave Operational Transformation is that it is able to compose the operations if they are compatible. This Composer is able to validate operations to make sure that they are compatible and compose them. One of the criteria for compatibility is that operations must span the same number of indexes. If operations a and b are composed into operation d , transformation function can take operations d, c and produce operation c'' that can be applied to peer's document and the document is in desired state.

This example is still simple. It does not point out some necessary parts of Operational transformation. The importance of preservation of operation c' was obvious from the diagram, but without the diagram it is hard to determine what operation should server or client preserve. In order to make it easier, metadata about the parent state of the operation is added. Every operation has an identifier of state in which it was executed. Google Wave uses a scalar version number, but [5] suggests to use a hash of the documents contents.

Having information about parent state, the server is able to determine that parent state of operation b is not in server history and it has to derive a new operation, that would take the document to state after application of operations a and c . This operation is c' . This deriving of operation c' is called bridging.[5]

With this new feature comes a new problem. In order to derive the operation c' server has to preserve also operation a . This could be a scalability problem in an environment with many clients editing the same document, since server has potentially preserve this data for every client. This is resolved by transferring all the responsibility to client and buffering the operations on client side. Client can send only operations that come from state in server's history and also he can send only one operation to the server and he has to wait for the acknowledgment of the operation. All following operations are preserved in buffer on client side. One of the advantages of buffering operations is that if these operations are compatible, they can be composed using the Composer. When the server acknowledges first operation, client can send next operation. This way client can predict server's path and send only operations that are on the server's path. Also it is much more complicated for the server to track every client, but it generally is much easier for client to track only one server. This way the server solves only one step diamond problem presented in the first example.

Daniel Spiewak in article Understanding and Applying Operational Transformation[5] describes some other improvements of Operational Transformation by Google Wave. For example, when client receives operation from server and has some operations in buffer, he can transform the whole buffer against this operation in order to accelerate the process on his site and also to ensure that operations in buffer have a parent state in servers history.

The main difference between Google control algorithm and the first dOPT algorithm is that dOPT doesn't use transformation function, it uses transformation matrix. If we have set of operations of cardinality m the matrix should be of size $m \times m$. The entries of the matrix correspond to all possible pairs of operations and contain functions which transform operations to other operations.

Another important part of the algorithm is the state vector. Timestamps for each client are handled by a vector timestamp, where a state vector s_i for a client C_i will have at position j , the number of operations known to have been executed by client C_j . [3]

According to this vector it is decided whether the received operation is executed or put in the request queue.

The difference between Google algorithm and Jupiter is not so significant. Actually, Jupiter also uses transformation function, it is called xform, although it is very similar to dOPT transformation matrix. Only difference is that Jupiter doesn't use the operation composition. According to [4] this algorithm assumes use of transport layer that delivers data in right order, such as TCP.

The last major algorithm implementing this technique is TIPS, which is modern and based on the admissibility-based transformation (ABT) framework. It is created on top of HTTP protocol. Main difference is that clients can join and leave session at any time and the algorithm is more reliable when network failures can be expected.

2.3 Differential Synchronization

Differential synchronization is the second technique used for real-time collaboration. It was fully described in 2009 by Neil Fraser in paper Differential Synchronization[2]. This section will describe and explain this technique, point out its advantages and disadvantages. As the main source of this information will be used paper mentioned above.

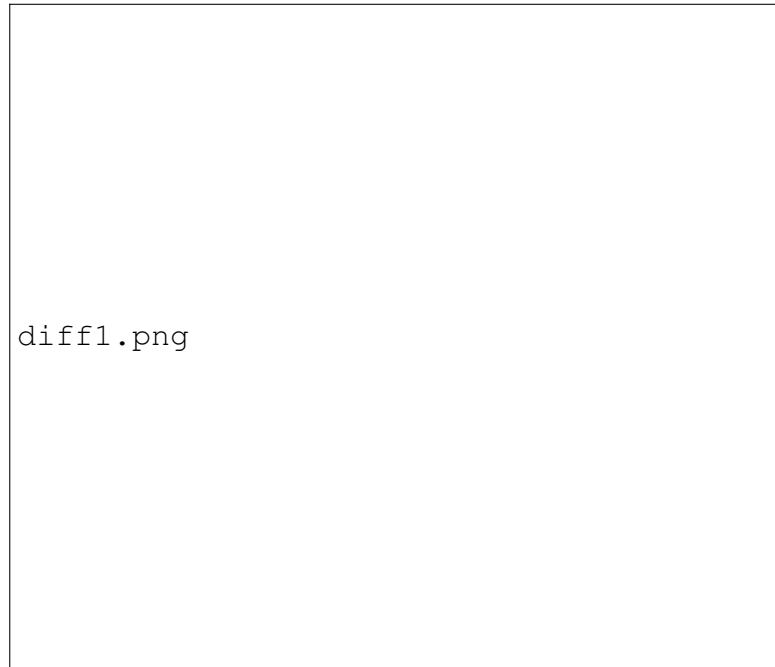
DS also replicates data on all sites of collaboration and identical code runs on both server and client, so it is symmetrical. The main building block is the presence of diff & patch algorithms. Diff algorithm is able to compute a difference between documents and save it in appropriate format for the patch algorithm, which is able to apply these changes on the other copy. The patch algorithm must be fuzzy, this means that the changes may be applied even if the document has changed in the meantime, because Differential Synchronization is an optimistic technique. The final version presented by Neil Fraser is suitable for unreliable networks and networks with high-latency. This property will be described later in this section.

The key feature of DS is that it is simple and well suited for use in both novel and existing state-based applications without requiring application redesign.[2] Differential Synchronization is able to process variety of formats, not only plain text. As long as there is a diff and patch algorithm available for the desired format, DS is able to use it. It is implemented and used in MobWrite.

Basic principle of this technique can be described using data flow diagram in Figure 2.4 originally presented in [2].

In the beginning Client text, Server text and Common Shadow are the same. Client text and Server text represent 2 sites of the collaboration. The goal is to keep them updated. Client and Server are allowed to make changes in their documents at any time.

Figure 2.4: Differential synchronization



After specified time interval a snapshot of Client text is taken and using a diff algorithm difference between snapshot and common shadow is computed. Common shadow represents the last common document state before any edits have been made. This diff returns changes that have been performed by Client. The snapshot is copied over to the common shadow and the changes are patched on the Server text. Now the process repeats symmetrically using the server text in order to apply changes made by server or other clients on the client text.

2.4 Commutative Replicated Data Types

3 Comparison of Techniques

3.1 Comparison Criteria

3.2 Comparison

4 Komodo

4.1 Overview

4.2 Komodo Requirements for Real-time Collaboration

4.3 Best Technique for Komodo

4.4 Java Implementation

5 Conclusion

Bibliography

- [1] ELLIS, C. A., AND GIBBS, S. J. Concurrency control in groupware systems. *SIGMOD Rec.* 18, 2 (June 1989), 399–407.
- [2] FRASER, N. Differential synchronization. In *Proceedings of the 9th ACM Symposium on Document Engineering* (New York, NY, USA, 2009), DocEng '09, ACM, pp. 13–20.
- [3] LEUNG, C. Operational transformation in cooperative software systems. *McGill Science Undergraduate Research Journal*, 8 (2013).
- [4] NICHOLS, D. A., CURTIS, P., DIXON, M., AND LAMPING, J. High-latency, low-bandwidth windowing in the jupiter collaboration system. In *Proceedings of the 8th Annual ACM Symposium on User Interface and Software Technology* (New York, NY, USA, 1995), UIST '95, ACM, pp. 111–120.
- [5] SPIEWAK, D. Understanding and applying operational transformation, 2015.
- [6] SUN, C., JIA, X., ZHANG, Y., YANG, Y., AND CHEN, D. Achieving convergence, causality preservation, and intention preservation in real-time cooperative editing systems. *ACM Trans. Comput.-Hum. Interact.* 5, 1 (Mar. 1998), 63–108.
- [7] VIDOT, N., CART, M., FERRIÉ, J., AND SULEIMAN, M. Copies convergence in a distributed real-time collaborative environment. 171–180.
- [8] WANG D., MAH A., L. S. Wave protocol google code, 2010.