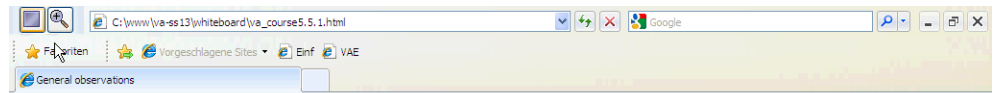# Script  generated by TTT

Title:  Distributed_Applications (10.06.2013)

Date:  Mon Jun 10 09:10:50 CEST 2013

Duration:  47:05 min

Pages:  15

---

Several requests to remote servers (e.g. RPC calls) may be bundled into a transaction.

```
begin-transaction
      callrpc (OP_1 , ...., )
      .....
      callrpc (OP_n , ...., )
end-transaction
```

A distributed transaction involves activities on multiple servers, i.e. within a transaction, services of several servers are utilized.

Transactions satisfy the *ACID* property: Atomicity, Consistency, Isolation, Durability.

1. *atomicity* : either all operations or no operation of the transaction is executed, i.e. the transaction is a success (commit) or else has no consequence (abort).

2. *durability* : the results of the transaction are persistent, even if afterwards a system failure occurs.

3. *isolation* : a not yet completed transaction does not influence other transactions; the effect of several concurrent transactions looks like as if they have been executed in sequence.

4. *consistency* : a transaction transfers the system from a consistent state to a new consistent state.

*Generated by Targeteam*

---

## Isolation

Isolation refers to the serializability of transactions. All involved servers are responsible for the serialization of distributed transactions. Example:

let U, T be distributed transactions accessing shared data on the two servers R and S.

if the transactions at server R are successfully executed in the sequence U before T, then the same commit sequence must apply to server S.

**Timestamp ordering**

**Locking**

**Optimistic concurrency control**

if conflicts are rare, optimistic concurrency control may be useful: no additional coordination necessary during transaction execution.

The check for access conflicts occurs when transactions are ready to "commit";

**Examples**

*Generated by Targeteam*

---

## Timestamp ordering

In a single server transaction, the server issues a unique timestamp to each transaction when it starts.

In a distributed transaction each server is able to issue globally unique timestamps.

for distributed transactions, the timestamp is the pair

(local timestamp, server-ID)

The local timestamp refers to the first server which issued the transaction timestamp.

Assume: $timestamp(trans) = t_{trans}$ and $timestamp(obj) = t_{obj}$

transaction trans accesses object obj

$$if (t_{trans} < t_{obj}) \ then \ abort(trans) \ else \ access \ obj;$$

*Generated by Targeteam*

Each server maintains locks for its own data items. Transaction trans requests lock (e.g. read, write lock) before access.

A transaction trans is well-formed if:

trans locks an object obj before accessing it.

trans does not lock an object obj which has already been locked by another transaction; except if the locks can coexist, e.g. two read locks.

prior to termination, trans removes all object locks.

A transaction is called a *2-phase* transaction if no additional locks are requested after the release of objects ("2-phase locking").

*Generated by Targeteam*

Isolation refers to the serializability of transactions. All involved servers are responsible for the serialization of distributed transactions. Example:

let U, T be distributed transactions accessing shared data on the two servers R and S.

if the transactions at server R are successfully executed in the sequence U before T, then the same commit sequence must apply to server S.

**Timestamp ordering**

**Locking**

**Optimistic concurrency control**

if conflicts are rare, optimistic concurrency control may be useful: no additional coordination necessary during transaction execution.

The check for access conflicts occurs when transactions are ready to "commit";

**Examples**

*Generated by Targeteam*

The following examples show the concurrency control approaches used by some current systems.

**Dropbox**

cloud service that provides file backup and enables users to share files and folders, accessing them from anywhere.

uses optimistic concurrency control; file granularity.

**Wikipedia**

creating and managing of wiki pages

uses optimistic concurrency control for editing.

**Google Docs**

cloud service providing web-based applications (word processor, spreadsheet and presentation) that allow users to collaborate by means of shared documents.

awareness based concurrency control: if several people edit the same document simultaneously, they will see each other's changes.

*Generated by Targeteam*

These aspects of distributed transactions may be realized by one of the following approaches. Let trans be a transaction.

**Intention list**

all object modifications performed by trans are entered into the intention list (log file).

When trans commits successfully, each server S performs all the modifications specified in $AL_S$ (trans) in order to update the local objects; the intention list $AL_S$ (trans) is deleted.

**New version**

When trans accesses the object obj, the server S creates the new version $obj_{trans}$ ; the new version is only visible to trans.

When trans commits successfully, $obj_{trans}$ becomes the new, commonly visible version of obj.

It trans aborts, $obj_{trans}$ is deleted.

*Generated by Targeteam*

This protocol supports the communication between all involved servers of the distributed transaction in order to jointly decide if the transaction should commit or abort.

We can distinguish between two phases

**Voting phase** : the servers submit their vote whether they are prepared to commit their part of the distributed transaction or they abort it.

**Completion phase** : it is decided whether the transaction can be successfully committed or it has to be aborted; all servers must carry out this decision.

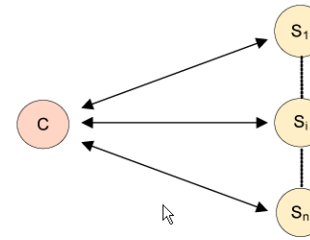**Steps of the two-phase commit protocol**

**Operations**

**Communication in the two-phase commit protocol**

**Problems**

---

One component (e.g. the client initiating the transaction or the first server in the transaction) becomes the **coordinator** for the commit process. In the following we assume, client C is the coordinator.



1. Coordinator C contacts all servers $S_i$ of the distributed transaction trans requesting their status for the commit (CanCommit?)
   - if server $S_k$ is not ready, i.e. it votes no, then the transaction part at $S_k$ is aborted;
   - $\exists$ i with $S_i$ is not ready

     then trans is aborted; the coordinator sends an abort message to all those servers who have voted with ready (i.e. yes).

2. $\forall$ i with $S_i$ is ready, i.e. commit transaction trans. Coordinator sends a commit message to all servers.

3. Servers send an acknowledgement to the coordinator.

---

The coordinator communicates with the participants to carry out the two-phase commit protocol by means of the following operations:

canCommit(trans) $\Rightarrow$ Yes/No: call from the coordinator to ask whether the participant can commit a transaction; participant replies with its vote.
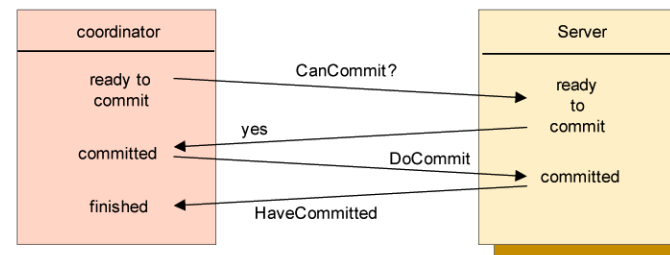
doCommit(trans): call from the coordinator to tell participant to commit its part of a transaction.

doAbort(trans): call from the coordinator to tell participant to abort its part of a transaction.

haveCommitted(trans, participant): call from participant to coordinator to confirm that it has committed the transaction.

getDecision(trans) $\Rightarrow$ Yes/No: call from participant to coordinator to ask for the decision on trans.

---

Number of messages: 4 * N messages for N servers.

During the 2PC process several failures may occur

one of servers crashes.

the coordinator crashes.

depending on their state, this may result in blocking situations, e.g. the coordinator waits for the commit acknowledge of a server, or a server waits for the final decision (commit or abort).

**Extended 2PC**

Three-Phase Commit protocol (3PC) is another approach to overcome blocking of servers until the crashed coordinator recovers.

*Generated by Targeteam*

---

```
Coordinator:
    multicast: ok to commit?
    collect replies
        all ok =>
                log commit to outcomes table
                wait until saved to persistent store
                send commit
        else  => send abort
    collect acknowledgements
    garbage collect data from outcomes table

After Failure:
    for each pending protocol in outcomes table
        send outcome (commit or abort)
        wait for acknowledgements
        garbage collect data from outcomes table
```

```
Server: first time message (CanCommit) received
    ok to commit =>
        save data to temp area (persistent store)
        reply ok
    commit =>
        make change permanent
        send acknowledgement
    abort  => delete temp area


message is a duplicate (recovering coordinator)
    send acknowledgement

After Failure:
    for each pending protocol
        contact coordinator to learn outcome
```

*Generated by Targeteam*

---

During the 2PC process several failures may occur

one of servers crashes.

the coordinator crashes.

depending on their state, this may result in blocking situations, e.g. the coordinator waits for the commit acknowledge of a server, or a server waits for the final decision (commit or abort).

**Extended 2PC**

Three-Phase Commit protocol (3PC) is another approach to overcome blocking of servers until the crashed coordinator recovers.

*Generated by Targeteam*

---

Distributed transactions are an important paradigm for designing reliable and fault tolerant distributed applications; particularly those distributed applications which access shared data concurrently.

**General observations**

**Isolation**

**Atomicity and persistence**

**Two-phase commit protocol (2PC)**

**Distributed Deadlock**

*Generated by Targeteam*