Introduction to STM

Part III: Operational Observations & Implementation Intricacies

Laurens Duijvesteijn June 2018

duijf.io · github.com/duijf · @_duijf

Preface

STM = Optimistic concurrency control

STM = Optimistic concurrency control

A TVar is a transactional variable

STM = Optimistic concurrency control

A TVar is a transactional variable

Operations on **TVar**s are called STM transactions

STM = Optimistic concurrency control

A TVar is a transactional variable

Operations on **TVar**s are called STM transactions

Run transactions with atomically; combine with orElse

Transactions get rerun on conflicts. We know because of a log

Transactions get rerun on conflicts. We know because of a log You can also call **retry** to rerun yourself

Transactions get rerun on conflicts. We know because of a log

You can also call **retry** to rerun yourself

The runtime reruns the transaction after its inputs change

Transactions get rerun on conflicts. We know because of a log

You can also call **retry** to rerun yourself

The runtime reruns the transaction after its inputs change

Use this for blocking operations

```
Recap (3/3)
TMVar a = TVar (Maybe a)
```

```
Recap (3/3)
TMVar a = TVar (Maybe a)
TQueue a = TQueue (TVar [a]) (TVar [a])
```

```
Recap (3/3)
```

```
TMVar a = TVar (Maybe a)
```

```
TQueue a = TQueue (TVar [a]) (TVar [a])
```

Want bounds? Add counts for read and write capacity!

```
Recap (3/3)
```

TMVar a = TVar (Maybe a)

TQueue a = TQueue (TVar [a]) (TVar [a])

Want bounds? Add counts for read and write capacity!

stm-containers works with a Hash Array Mapped Trie to avoid contention on structural updates

Today (by request)

Operational semantics of STM

Implementation notes

Again: This is all work from other people (Composable Memory Transactions, 2006)

Operational Semantics

Can we be precice about what code is doing?

Code is all well and good, words are nice too

But sometimes we need to specify things precicely

Can we be precice about what code is doing?

Code is all well and good, words are nice too

But sometimes we need to specify things precicely

Math and logic can help communicate clearly what's going on

A model can help us reason about what the code is doing

Where do models come from?

The authors create them! They come up with a language to communicate/specify what the code is doing.

It is a communication tool.

What does a model mean for me?

In the end, it is just that: a model

The implementation can differ from the model.

What does a model mean for me?

In the end, it is just that: a model

The implementation can differ from the model.

Or: the model doesn't cover everything in the implementation

What does a model mean for me?

In the end, it is just that: a model

The implementation can differ from the model.

Or: the model doesn't cover everything in the implementation Some people try very hard to prove an implementation corresponds to a model

Confession time!

For me, this is the first time I looked into these things

The math does look scary; I'm not sure I got everything right

Confession time!

For me, this is the first time I looked into these things

The math does look scary; I'm not sure I got everything right

This is a "best effort" attempt; hopefully it is still interesting and mostly correct. Otherwise, I'll stand corrected

Why do you even talk about this, if you're not an expert?

Why do you even talk about this, if you're not an expert?

Precicely because I'm not an expert.

Why do you even talk about this, if you're not an expert?

Precicely because I'm not an expert.

This stuff can be intimidating (it was to me!)

If this helps show you that **you too** can take a naive initial step at understanding intimidating stuff; that's a win :)

With that long disclaimer..

And apologies to the original authors;

With that long disclaimer...

And apologies to the original authors; other people upon whose toes I'm stepping;

With that long disclaimer..

And apologies to the original authors; other people upon whose toes I'm stepping; and possibly this audience;

With that long disclaimer..

And apologies to the original authors; other people upon whose toes I'm stepping; and possibly this audience; let's get to it.

Different kinds of semantics

There's multiple ways of specifying programming languages
The one we're looking at here is called "small-step
operational semantics"

Different kinds of semantics

There's multiple ways of specifying programming languages

The one we're looking at here is called "small-step operational semantics"

It looks like this...

```
P: \Theta \xrightarrow{a} O: \Theta'
                                                                                                           I/O transitions
                                                     \mathbb{P}[\operatorname{putChar} c]; \Theta \xrightarrow{!c} \mathbb{P}[\operatorname{return} ()]: \Theta
                                                                                                                                                                                                                      (PUTC)
                                                      \begin{array}{ccc} \mathbb{P}[\mathsf{getChar}]; \; \Theta & \stackrel{?c}{\longrightarrow} & \mathbb{P}[\mathsf{return}\; c]; \; \Theta \\ \mathbb{P}[\mathsf{forkI0}\; M]; \; \Theta & \longrightarrow & (\mathbb{P}[\mathsf{return}\; t] \mid M_t); \; \Theta & t \not\in \mathbb{P}, \Theta, M \end{array} \right. 
                                                                                                        \frac{M \to N}{\mathbb{P}[M] \colon \Theta \to \mathbb{P}[N] \colon \Theta} (ADMIN)
       \frac{M; \Theta, \{\} \stackrel{*}{\Rightarrow} \text{ return } N; \Theta', \Delta'}{\mathbb{P}[\text{atomic } M]: \Theta \rightarrow \mathbb{P}[\text{return } N]: \Theta'} (ARET) \qquad \frac{M; \Theta, \{\} \stackrel{*}{\Rightarrow} \text{ throw } N; \Theta', \Delta'}{\mathbb{P}[\text{atomic } M]: \Theta \rightarrow \mathbb{P}[\text{throw } N]: \Theta \cup \Delta'} (ATHROW)
                                                                                                     Administrative transitions
                                                                                                                                                                          M \rightarrow N
                                                                                   M \rightarrow V \text{ if } V[M] = V \text{ and } M \not\equiv V \quad (EVAL)
    return N >>= M \longrightarrow MN
                                                                                                    (BIND)
                                                                                                                                              \operatorname{catch} (\operatorname{return} M) N \rightarrow \operatorname{return} M (\operatorname{CATCH1})
                                                                                                                                            \operatorname{catch} (\operatorname{throw} M) N \to NM
       throw N >>= M \rightarrow \text{throw } N \pmod{THROW}
                                                                                                                                                                                                                                                                             (CATCH2)
              retrv >>= M \rightarrow retrv
                                                                                                     (RETRY)'
                                                                                                                                                               \operatorname{catch} \operatorname{retrv} N \longrightarrow \operatorname{retrv}
                                                                                                                                                                                                                                                                             (CATCH3)
                                                                                               STM transitions M: \Theta, \Delta \Rightarrow N: \Theta', \Delta'
                       \begin{array}{ccc} \mathbb{E}[\operatorname{readTVar}\,r];\;\Theta,\Delta &\Rightarrow & \mathbb{E}[\operatorname{return}\,\Theta(r)];\;\Theta,\Delta & \text{if } r \in dom(\Theta) & (READ) \\ \mathbb{E}[\operatorname{writeTVar}\,rM];\;\Theta,\Delta &\Rightarrow & \mathbb{E}[\operatorname{return}\,r];\;\Theta[r \mapsto M],\Delta & \text{if } r \in dom(\Theta) & (WRITE) \\ \mathbb{E}[\operatorname{newTVar}\,M];\;\Theta,\Delta &\Rightarrow & \mathbb{E}[\operatorname{return}\,r];\;\Theta[r \mapsto M],\Delta[r \mapsto M] & r \not\in dom(\Theta) & (NEW) \end{array}
           \frac{M \to N}{\mathbb{E}[M]: \Theta, \Delta \Rightarrow \mathbb{E}[N]: \Theta, \Delta} (AADMIN) \qquad \frac{M_1; \Theta, \Delta \Rightarrow \mathsf{return} \ N; \Theta', \Delta'}{\mathbb{E}[M_1: \mathsf{orElse}^t \ M_2]: \Theta, \Delta \Rightarrow \mathbb{E}[\mathsf{return} \ N]: \Theta', \Delta'} (OR1)
\frac{M_1; \Theta, \Delta \stackrel{\Rightarrow}{\Rightarrow} \text{throw } N; \Theta', \Delta'}{\mathbb{E}[M_1 \text{ 'orElse' } M_2]; \Theta, \Delta \Rightarrow \mathbb{E}[\text{throw } N]; \Theta', \Delta'} (OR2) \qquad \frac{M_1; \Theta, \Delta \stackrel{\Rightarrow}{\Rightarrow} \text{retry; } \Theta', \Delta'}{\mathbb{E}[M_1 \text{ 'orElse' } M_2]; \Theta, \Delta \Rightarrow \mathbb{E}[M_2]; \Theta, \Delta} (OR3)
```

Figure 4: Operational semantics of STM Haskell

First off...

This is all about single threaded Haskell (the 2006 paper)

First off...

This is all about single threaded Haskell (the 2006 paper)

The runtime: repeatedly picks a thread, then executes an IO transition. It ends up in a new state.

So: a transition is a way to move from one state to another

First off...

This is all about single threaded Haskell (the 2006 paper)

The runtime: repeatedly picks a thread, then executes an IO transition. It ends up in a new state.

So: a transition is a way to move from one state to another

There are **two** types of transitions. \Rightarrow is for STM; \rightarrow is for IO.

We also saw some stuff with bars

These are called rules of inference.

We also saw some stuff with bars

These are called rules of inference. They work like this:

___premises __conclusion

We also saw some stuff with bars

These are called rules of inference. They work like this:

premises conclusion

In logic, you can have multiple premises. Today, we only look at rules with one.

$$M;\Theta,\{\}\Rightarrow^* pure\ N;\Theta',\Delta'$$
 $\mathbb{P}[atomic\ M];\Theta\to\mathbb{P}[pure\ N];\Theta'$
 (ret)

$$M; \Theta, \{\} \Rightarrow^* \mathsf{throw} \ N; \Theta', \Delta'$$
 $\mathbb{P}[\mathsf{atomic} \ M]; \Theta \to \mathbb{P}[\mathsf{throw} \ N]; \Theta \cup \Delta'$
 (throw)

So far, we've learned/specified

The heap can only be modified by an **atomic** block (from STM)

So far, we've learned/specified

The heap can only be modified by an **atomic** block (from STM)

Exceptions propagate from STM, and leave allocation effects (new TVars, as we'll see soon) but no other modifications.

Next up is reading/writing/creating TVars

$$\mathbb{E}[\mathsf{readTVar}\ r]; \Theta, \Delta \Rightarrow \mathbb{E}[\mathsf{pure}\ \Theta(r)]; \Theta, \Delta$$

$$\mathsf{if}\ r \in dom(\Theta)$$

$$(\mathsf{read})$$

$$\mathbb{E}[\mathsf{writeTVar}\ r\ M];\Theta,\Delta\Rightarrow\mathbb{E}[\mathsf{pure}\ (\)];\Theta[r\mapsto M],\Delta$$
 if $r\in dom(\Theta)$ (write)

$$\mathbb{E}[\mathsf{newTVar}\ M]; \Theta, \Delta \Rightarrow \mathbb{E}[\mathsf{pure}\ r]; \Theta[r \mapsto M], \Delta[r \mapsto M]$$
 if $r \not\in dom(\Theta)$ (new)

Another pitstop

We've seen creation, reading and writing

Also: were those allocation effects came from

Next stop is orElse, which is more involved

orElse has three cases

- 1. The first transaction succeeds
- 2. The first transaction throws
- 3. The first transaction retries

$$\frac{M_1; \Theta, \Delta \Rightarrow^* \text{pure } N; \Theta', \Delta'}{\mathbb{E}[M_1 \text{ `orElse`} M_2]; \Theta, \Delta \Rightarrow \mathbb{E}[\text{pure } N]; \Theta', \Delta'}$$
(or1)

$$\begin{array}{c} \textit{M}_1;\Theta,\Delta\Rightarrow^*\mathsf{throw}\ \textit{N};\Theta',\Delta'\\ \hline \mathbb{E}[\textit{M}_1\ \ \mathsf{`orElse}\ \ \textit{M}_2];\Theta,\Delta\Rightarrow\mathbb{E}[\mathsf{throw}\ \textit{N}];\Theta',\Delta'\\ \hline \\ (\mathsf{or2}) \end{array}$$

Why is there a Θ' there, and not a Θ ?

Why is there a Θ' there, and not a Θ ?

This means that the heap modifications are still there. Is that a problem?

Why is there a Θ' there, and not a Θ ?

This means that the heap modifications are still there. Is that a problem?

Remember the rule from before (in IO instead of STM)?

Why is there a Θ' there, and not a Θ ?

This means that the heap modifications are still there. Is that a problem?

Remember the rule from before (in **IO** instead of **STM**)?That addresses this.

$$M_1; \Theta, \Delta \Rightarrow^* \text{retry}; \Theta', \Delta'$$

$$\mathbb{E}[M_1 \text{ `orElse` } M_2]; \Theta, \Delta \Rightarrow \mathbb{E}[M_2]; \Theta, \Delta$$

$$(\text{or3})$$

There's something missing, though

There's no explicit rule for what happens when a **retry** is reached, while in **IO**.

There's something missing, though

There's no explicit rule for what happens when a **retry** is reached, while in **IO**.

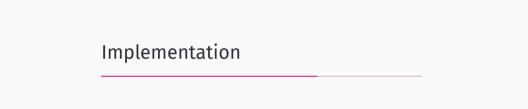
Because the system cannot make any progress when this happens. It blocks.

There's something missing, though

There's no explicit rule for what happens when a **retry** is reached, while in **IO**.

Because the system cannot make any progress when this happens. It blocks.

Other threads can make progress, and if any **TVar** of the transaction changes, the thread can proceed.



So I mentioned this log a couple of times..

Its purpose is to find conflicts before a transaction commits Let's see how it works

Each thread has a log

It gets created as soon as a transaction starts (e.g. the atomic block is encountered)

Each log contains:

- 1. A parent pointer (may be null)
- 2. A number of log entries

Log entries

One exists for each **TVar** that's touched by the transaction.

These have three things:

- 1. The TVar
- 2. Its value when the transaction started
- 3. Its current value

How do we know there's conflicts?

For a transaction, we check all entries in the log.

There's a conflict if one of the **TVar**s changed during the transaction.

How do we know there's conflicts?

For a transaction, we check all entries in the log.

There's a conflict if one of the **TVar**s changed during the transaction.

(Checking this is linear in the amount of **TVar**s accessed)

Aborting transactions

Aborting a transaction simply means discarding its log

Aborting transactions

Aborting a transaction simply means discarding its log

Writes happen to the log, in the "new value once transaction commits" slot, so they don't have any effects.

Aborting transactions

Aborting a transaction simply means discarding its log

Writes happen to the log, in the "new value once transaction commits" slot, so they don't have any effects.

Reads also happen from the log, so a transaction can see the values it has written. This also helps with A-B-A problems.

Committing transactions

Loop over the transaction log and, for each entry, write the new value to the **TVar**.

In the single threaded runtime, this is enough. Only one thread can commit at the same time.

Multicore would need some sort of synchronization: or you'll have trouble when two threads commit to the same **TVar**s in different orders

Nontermination guarantees

There's cases where transactions can loop indefinitely because of the values they read.

Nontermination guarantees

There's cases where transactions can loop indefinitely because of the values they read.

It's obvious after you see it, but I never thought of it.

Nontermination guarantees

There's cases where transactions can loop indefinitely because of the values they read.

It's obvious after you see it, but I never thought of it.

Fix: before switching to a thread in a transaction, run the check routine. (remember: single threaded runtime)

Avoiding busy wait

We mentioned: there's no busy wait because transactions are only reran when their inputs change.

Avoiding busy wait

We mentioned: there's no busy wait because transactions are only reran when their inputs change.

But: how does the runtime know which threads to rerun?

Avoiding busy wait

We mentioned: there's no busy wait because transactions are only reran when their inputs change.

But: how does the runtime know which threads to rerun?

Each **TVar** has a list of blocking transactions! When a thread commits to a **TVar**, it wakes up another.

The presence of **orElse** means we have to have some form of transaction nesting

The presence of **orElse** means we have to have some form of transaction nesting (to isolate excecution of alternatives).

The presence of **orElse** means we have to have some form of transaction nesting (to isolate excecution of alternatives).

For this, we create a nested log. (Remember the parent pointer?)

The presence of **orElse** means we have to have some form of transaction nesting (to isolate excecution of alternatives).

For this, we create a nested log. (Remember the parent pointer?)

If a nested transaction finishes, we validate its log and that of the parent. Both good? Then merge the two logs.

What **TVar**s should we wait on?

What **TVar**s should we wait on? A: The union of those accessed in both transactions!

What **TVar**s should we wait on? A: The union of those accessed in both transactions!

This is so all the **TVar**s which have been accessed can get the current thread ID on their list of waiting threads.

What **TVar**s should we wait on? A: The union of those accessed in both transactions!

This is so all the **TVar**s which have been accessed can get the current thread ID on their list of waiting threads.

Observe: this log is only used for installing these waiting entries, or its discarded.

We've seen: Operational semantics of STM, and some implementation notes of the single threaded version.

We've seen: Operational semantics of STM, and some implementation notes of the single threaded version.

That's not how the real thing works, but it is a good didactic model IMO.

Hopefully this has been interesting to you!

The slides are all going to be online: duijf/stm-course

Hopefully this has been interesting to you!

The slides are all going to be online: duijf/stm-course

Feedback: feel free to tell me or write on a GitHub ticket!

Fin.