

# Implementation of a Linear Session Type System

Second assignment for the course in Languages for  
Concurrency and Distribution, A.Y. 2022/2023

Christian Micheletti  
July 9, 2024



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

- 1 Structure
- 2 Language
  - Syntax
  - Runtime
- 3 Types
  - Basic types
  - Naive Implementation
- 4 Algorithmic Type Checking
- 5 Other features
  - Recursive Types
  - Tuple Syntax
  - Replicated Behaviour

The paper *Fundamentals of Session Types* presents a  $\pi$ -calculus equipped with a (linear) type system to check against errors that the system can incur into.

Such language features linear types, which allow checking whether a given channel is used by exactly one process.

This  $\pi$ -calculus' features are incrementally presented through five sections:

- language **syntax and semantics**, and **basic types and rules**, in sections 2 and 3;
- **recursive types**, needed to write useful unrestricted types, in section 4;
- **unbounded computations**, in section 5, that enable infinite behaviours;
- **branching and selection**, in section 6 (not implemented).

The implementation process follows this very same organization, and will be presented accordingly.

The software is written in Haskell and features:

- a **code interpreter**, that executes the code in Concurrent Haskell;
- two **type checkers**:
  - one following the rules with context split;
  - one implementing the algorithmic rules presented in **section 8**;
- a type inferrer.

The paper has several examples of the expected behaviour of the type checker, that are implemented as unit tests.

Section 2's original syntax

```
P ::=  $\bar{x} v . P$   
       $x(x) . P$   
       $P \mid P$   
      if  $v$  then  $P$  else  $P$   
      0  
       $(\nu x x) P$ 
```

```
v ::= x  
      true  
      false
```

Implementation

```
P ::= x << v . P  
      | x >> x . P  
      | P | P  
      | if v then P else P  
      | 0  
      | x << x . P  
      | { P }
```

```
v ::= x  
      | true  
      | false
```

After parsing, the program is precomputed lifting all bindings over all parallel compositions, as expressed in the following **structural congruence**:

$$(\nu xy)P|Q \equiv (\nu xy)(P|Q)$$

The code interpreter prints debug information about the program executed, along with a timestamp and a description of what the process did at that time.

## Output format

```
[TIMESTAMP | ThreadId THREAD_ID]: MESSAGE
```

Program behaviour is defined according to the **operational semantics** presented in section 2.



The language runtime is implemented in concurrent Haskell, using:

- the IO monad, along with a local state that maps variables to channels and literals, to model threads;
- and the MVar type to model channels.

Each thread is created with the

```
forkIO :: IO () -> IO ThreadId
```

function, that simply creates a new thread, returning its id.

Plain MVars can await for value insertion with the function

```
takeMVar :: MVar a -> IO a
```

but the dual operation

```
putMVar :: MVar a -> a -> IO ()
```

does not await the variable to be ready to accept a new value.

Channels are represented as tuples of `MVar`, meaning `(value, idle)`:

- `value :: MVar v`, containing the passed value, which needs to be evaluated;
- `idle :: MVar ()`, which has a `()` value in it when the channel is ready to receive.

Then, the `Channel` datatype further distinguishes whether it is a read end or a write end. Read ends and write ends share the same `MVars`.

The implementation is as follows:

```
receive :: ChannelEnd -> IO Val
receive (ReadEnd var idle) = do
    putMVar idle ()      -- signals the other end that it
                        -- is ready to receive

    takeMVar var
receive _ = throw $ userError "..."
```

  

```
send :: ChannelEnd -> Val -> IO ()
send (WriteEnd var idle) val = do
    takeMVar idle        -- waits until the other end is
                        -- ready to receive

    putMVar var val
send _ _ = throw $ userError "..."
```

## Inaction

The process 0 just prints STOP and ends the thread.

## Branching

The process if  $v$  then  $P_1$  else  $P_2$  prints two debug messages, both starting with BRANCHING:

- the guard before evaluation;
- and then, after evaluation in the local state.

After that, the run continues as the appropriate process.

## Binding

The process  $(\nu xy).P$  just prints BINDING followed by the two bounded variables. Then, it creates two MVars, one for the value and one for the lock and associates in the local state both variables with the respective ends.

The run proceeds in the same thread as prescribed by  $P$ .

## Fork

The process  $P_1|P_2$  prints FORK followed by the two new processes ids.

To prevent a concurrent program to end before all forked threads terminate, this code constructs two MVars that are notified when the two new threads finish. This process doesn't do anything more than awaiting for both threads to finish.

## Sending

The process  $\bar{x}v.P$  prints SENDING followed by the value:

- before evaluation;
- and after evaluation.

Then the process sends the value over the channel  $x$  and proceeds as prescribed by  $P$ .

## Receiving

The process  $x(v).P$  prints RECEIVING followed by:

- the newly bound variable name;
- and the value received.

Then the process proceeds as prescribed by  $P$ .

Implementation of qualifier, pretypes and types are fully compliant to section 3's original syntax:

$$q ::= \text{lin}$$
$$\text{un}$$
$$p ::= ?T.T$$
$$!T.T$$
$$T ::= \text{bool}$$
$$\text{end}$$
$$qp$$
$$\Gamma ::= \emptyset$$
$$\Gamma, x : T$$

Contexts are implemented as hash maps of types



Duality is partially defined as follows:

```
dualType :: SpiType -> SpiType
dualType End = End
dualType Boolean = error "..."
dualType (Qualified q (Receiving t1 t2)) =
    Qualified q (Sending t1 (dualType t2))
dualType (Qualified q (Sending t1 t2)) =
    Qualified q (Receiving t1 (dualType t2))
dualType (Recursive a p) = Recursive a (dualType p)
dualType (TypeVar x) = TypeVar
```

Contexts support the following operations:

- (nondeterministic) **context split**  $\Gamma = \Gamma_1 \circ \Gamma_2$ : the function `ndsplitt :: Context -> [(Context, Context)]` creates all possible combinations of dividing linear variables, maintaining unrestricted variables;
- **update**  $\Gamma + (x : T)$ : the function `update k t` inserts  $k : t$  in the context only if the variable  $k$  was not present, or it was yet defined unrestricted with type  $t$ ;
- **override**  $\Gamma, (x : T)$ : this represents newly bounded variables, possibly shadowing preexisting definitions.

Sequent calculus rules and operations over contexts are modeled as instances of the context transition **monad**:

```
newtype CT a = CT
  (Context -> TypeErrorBundle TypeError (a, Context))
```

```
instance Monad CT where
  return :: a -> CT a
  (>>=) :: CT a -> (a -> CT b) -> CT b
```

Rules can be composed and propagate context side effects, and can yield an output

```
class TypeCheck a where
  type Output a
  check :: a -> CT (Output a)
```

## Unrestricted requirement

$\text{un}(\Gamma)$  holds when all entries in the context are unrestricted:

```
unGamma :: CT ()
unGamma = CT (\context -> if all unrestricted context
  then Right ((), context)
  else Left "Error message...")
```

## Context update

$\Gamma + (x : T)$  throws an error if the conditions aren't met:

```
update :: String -> SpiType -> CT ()
update k t = do
  may <- liftPure (M.lookup k)
  case may of
    Just found -> unless (unrestricted t && found ≈ t)
      (throwError "Error message..")
    Nothing     -> sideEffect (M.insert k t)
```

## Variable extraction

Context transitions can even return useful values:

```
extract :: String -> CT SpiType
extract k = do
  t <- get k
  unless (unrestricted t) (delete k)
  return t
```

This is used to optimize some context splits in the [T-REC] and [T-SEND] rules

The paper shows two sets of rules to type check a process:

- the rules based on **context split**;
  - naive approach, requires to check for all possible context splits for parallel composition (other rules are optimizable);
  - can be done without using the Output of the algorithm;
- the rules presented in section 8 (Algorithmic Type Checking):
  - more efficient, doesn't rely on context split;
  - requires to track all used variables along the program and return them as output.

Context split rules are modelled as  $CT()$ , i.e. transitions that yield no output.

One could think of an intuitive way to parallelize the  $[T-PAR]$ :

$[T-PAR]$

$$\frac{\Gamma_1 \vdash P_1 \quad \Gamma_2 \vdash P_2}{\Gamma_1 \circ \Gamma_2 \vdash P_1 | P_2} [T-PAR]$$



The algorithm checks all possible splits of  $\Gamma_1 \circ \Gamma_2$

```
check (Par p1 p2) = do
  splits <- liftPure ndsplit
  -- Compute all possible splits
  runs <- return () -< (candidate <$> splits)
  liftEither $ foldChoice runs
    {- `using` parList rdeepseq -}
  where
    candidate (c1, c2) = (return () -<
      [ c1 |> check p1
        , c2 |> check p2
      ] {- <del></del> `using` parList rdeepseq -}
    ) >- return ()
```

It seems natural to desire to parallelize this code in the points with comments...

... but this would imply loss in performance. Consider the following program:

## Example: assets/well-formed-ill-typed/multiple.spi

```
a1 << a2: lin?bool.end .
b1 << b2: lin?bool.end .
c1 << c2: lin?bool.end .
d1 << d2: lin?bool.end .
e1 << e2: lin?bool.end .
f1 << f2: lin?bool.end .
x << y: rec x. ?bool.x .
  x1 << y1: lin?bool.lin!bool.end .
  x2 << y2: lin?bool.lin!bool.end .
    { x << true . y >> z . if z then 0 else 0
      | y >> z . if z
        then x << false . 0
        else 0
      | x1 << true . x1 >> n . y2 >> n . y2 << n . 0
      | y1 >> n . y1 << false . x2 << false . x2 >> n . 0
      | a1 << true . b1 << true . c1 << true
        . d1 << true . e1 << true . f1 << true . 0
      | a2 >> e . b2 >> e . c2 >> e . d2 >> e . e2 >> e .
        0
    }
  }
```

## Commands to run

```
stack build  
stack exec lcd-assignment-exe -- +RTS -s -ls -N2 -RTS \  
    nd assets/well-formed-ill-typed/multiple.spi
```

Running this program with one CPU, it takes roughly 7 seconds to terminate.

Running this program with more CPUs makes it even worse!

Other rules are optimized to not resort to context split.

## Example: T-IN

$$\frac{\Gamma_1 \vdash x : q?T.U \quad (\Gamma_2 + x : U), y : T \vdash P}{\Gamma_1 \circ \Gamma_2 \vdash x(y).P} \text{ [T-IN]}$$

In this case the rule holds following those observations:

- if  $x$  is not present in  $\Gamma_1 \circ \Gamma_2$  the left assumption can never be verified;
- if  $\Gamma_1 \circ \Gamma_2$  contains the claim  $x : \text{un}T.U$ , then both  $\Gamma_1, \Gamma_2$  contain such claim;
- if  $\Gamma_1 \circ \Gamma_2$  contains the claim  $x : \text{lin}T.U$ , then  $\Gamma_1$  must contain that claim and  $\Gamma_2$  must not.

Notice that in order to type  $\Gamma_1 \vdash x : q?T.U$ , as for the rule [T-VAR],  $\Gamma_1$  must not contain any linear claim.

Hence,  $\Gamma_2$  must contain all linear claims in  $\Gamma_1 \circ \Gamma_2$ , of course except for  $x$  if it was linear.

Hence the algorithm for this rule is as follows:

## Example: T-IN

```
check (Rec x y p) = do
  xType <- extract x -- If not present, the monadic
    -- bind will make the whole rule fail.
    -- The function extract will preserve
    -- any unbounded channel
  (t, u) = case xType of
    Qualified _ (Receiving t u) ->
      return (t, u)
    _ -> throwError "..."
```

update x u  
replace y t  
check p

The problem with the parallel composition rule is that we can't know a priori which parallel branch uses a variable.

However, we can easily avoid this problem by keeping track of **variables** used in **subject position**, such as channel send or receive.

For processes, output context  $\Gamma_2$  contains all free variables “consumed” by process  $P$ ,  $L$  contains all linear free variable used in process  $P$

$$\Gamma_1 \vdash P : \Gamma_2; L$$

$$\Gamma_1 \vdash x : T; \Gamma_2$$

For variables, the output is a context without the variable, if it was linear

Then **context difference**  $\Gamma \div L$  removes all variables in  $L$  from  $\Gamma$  and yields an error if there was a linear variable in  $L$ .

Then, parallel composition rule is easily defined:

$$\frac{\Gamma_1 \vdash P : \Gamma_2; L_1 \quad \Gamma_2 \div L_1 \vdash Q : \Gamma_3; L_2}{\Gamma_1 \vdash P|Q : \Gamma_3; L_2} \text{ [A-PAR]}$$

The implementation is straight forward:

```
check (Par p1 p2) = do
  l1 <- check p1
  contextDiff l1
  check p2
```



Other rules have similar implementation. The optimizations for direct context split check of [T-IN] and [T-OUT] are related.

Rule [A-IF] requires that the outputs for both branches are the same:

$$\frac{\Gamma_1 \vdash v : q \text{ bool}; \Gamma_2 \quad \Gamma_2 \vdash P : \Gamma_3; L \quad \Gamma_2 \vdash Q : \Gamma_3; L}{\Gamma_1 \vdash \text{if } v \text{ then } P \text{ else } Q : \Gamma_3; L} \text{ [A-IF]}$$

Other features include:

- **recursive types** (described in section 4);
- **tuple syntax** (described in section 4);
- **replicated behaviours** (described in section 5);
- **branching and selection** (described in section 6 but not implemented).

Context update doesn't allow [T-IN] and [T-OUT] rules (and algorithmic respectives) to type any unbounded type:

## Rule T-IN

$$\frac{\Gamma_1 \vdash x : q?T.U \quad (\Gamma_2 + x : U), y : T \vdash P}{\Gamma_1 \circ \Gamma_2 \vdash x(y).P} \text{ [T-IN]}$$

It requires the  $U$  type to be the same as  $un?T.U$ . Same holds for [T-OUT].

We would like to consider the solution of an equation like

$$x = \text{un?}T.x$$

In order to do this, the paper introduces two new syntax constructs:

- **type variables;**
- a type-level binder  $\mu$  that define types as **infinite regular trees**.

For example, a type that sends a boolean type all over can be described just as follows:

$$\mu x. \text{un!} \text{bool}. x \qquad \text{rec } x . \text{ un ! bool } . x$$

The paper omits the co-inductive definition, the code implements this equivalence as a **bisimulation**

```
class (Eq a, Ord a) => Bisimulation a where
  behave :: a -> (a, a)

bisim :: (Bisimulation a) => S.Set (a, a) -> a -> a -> Bool
bisim _ a b | a == b = True
bisim rel a b | S.member (a, b) rel = True
bisim rel a b | S.member (b, a) rel = True
bisim rel a b =
  let (a', oa) = behave a
      (b', ob) = behave b
  in (oa == ob && bisim (S.insert (a, b) rel) a' b')

(≈) :: (Bisimulation a) => a -> a -> Bool
(≈) = bisim S.empty
```

Recursive types  $\mu a.T$  implement the `behave` function required, as specified in the paper, to never consider the type as it is but another type in the same equivalence class, namely  $T[\mu a.T/a]$

```
instance Bisimulation SpiType where
  behave :: SpiType -> (SpiType, SpiType)
  behave (Qualified q (Sending v t)) =
    (t, Qualified q (Sending v End))
  behave (Qualified q (Receiving v t)) =
    (t, Qualified q (Receiving v End))
  behave t@(Recursive x t') = behave (subType x t t')
  behave t = (t, t)
```

The paper presents a syntactic sugar for tuple passing:

$\overline{x_1}\langle u, v \rangle.P$  abbreviates  $(\nu y_1 y_2) \overline{x_1} y_2. \overline{y_1} u. \overline{y_1} v. P$

$x_2(w, t).P$  abbreviates  $x_2(z).z(w).z(t).P$

Hidden channels  $y_1, y_2$  need to be typed. There are two ways to achieve this:

- type annotation on the binding  $(\nu x_1 x_2)$  (long and tedious);
- type inference (not described in the paper, hence details omitted).

Section 5 modifies process syntax annotating receiving channels with `lin` and `un` modifiers:

- if a channel is annotated with `lin` then it is consumed;
- if a channel is annotated with `un` then it is replicated.

Reduction rules (and runtime implementation) are updated accordingly:

$$(\nu xy)\bar{x}v.P|\text{lin}y(z).Q \rightarrow (\nu xy)P|Q[v/z] \quad [\text{R-LINCOM}]$$

$$(\nu xy)\bar{x}v.P|\text{un}y(z).Q \rightarrow (\nu xy)(P|Q[v/z]|\text{un}y(z).Q) \quad [\text{R-UNCOM}]$$

(In practice, only the `un` modifier is introduced)



Types are updated following a third invariant:

**inv. (iii)** Unrestricted input processes may not contain free linear variables.

That implies:

- there is no replication of linear variables, maintaining their transient behaviour;
- there could be *bounded* linear variables.

Hence the new rules for receiving channels are:

## Rule T-IN

$$\frac{q(\Gamma_1 \circ \Gamma_2) \quad \Gamma_1 \vdash x : q?T.U \quad (\Gamma_2 + x : U), y : T \vdash P}{\Gamma_1 \circ \Gamma_2 \vdash qx(y).P} \text{ [T-IN]}$$

## Rule A-IN

$$\frac{\begin{array}{l} \Gamma_1 \vdash x : q_2?T.U; \Gamma_2 \\ (\Gamma_2, y : T) + x : U \vdash P : \Gamma_3; L \end{array} \quad q_1 = \text{un} \implies L \setminus \{y\} = \emptyset}{\Gamma_1 \vdash q_1x(y).P : \Gamma_3 \div \{y\}; L \setminus \{y\} \cup (\text{if } q_2 = \text{lin then } \{x\} \text{ else } \emptyset)} \text{ [A-IN]}$$

To run the program the command is:

```
stack run -- [nd] path/to/file.spi
```

Specify the `nd` flag to use the context split type checker.