Dependent Session Types for Verified Concurrent Programming

ANONYMOUS AUTHOR(S)

We present TLL_C which extends the Two-Level Linear dependent type theory (TLL) with session type based concurrency. Equipped with Martin-Löf style dependency, the session types of TLL_C allow protocols to specify the properties of communicated messages. When used in conjunction with the dependent type machinery already present in TLL, dependent session types facilitate the a form of relational verification by relating concurrent programs with their idealized sequential counterparts. Correctness properties proven for sequential programs can now be easily lifted to their corresponding concurrent programs. Session types now become a powerful tool for intrinsically verifying the correctness of data structures such as queues and concurrent algorithms such as map-reduce. To extend TLL with session types, we develop a novel formulation of intuitionistic session type which we believe to be widely applicable for integrating session types into other type systems beyond the context of TLL_C . We study the meta-theory of our language, proving its soundness as both a term calculus and a process calculus. All reported results are formalized in Rocq. A prototype compiler which compiles TLL_C programs into concurrent C code is implemented and freely available.

Additional Key Words and Phrases: dependent types, linear types, session types, concurrency

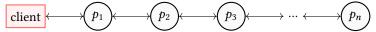
1 Introduction

Session types [23] are an effective typing discipline for coordinating concurrent computation. Through type checking, processes are forced to adhere to communication protocols and maintain synchronization. This allows session type systems to statically rule out runtime bugs for concurrent programs similarly to how standard type systems rule out bugs for sequential programs. While (simple) session type systems guarantee concurrent programs do not crash catastrophically, it remains difficult to write concurrent programs which are semantically correct.

Consider the Pfenning-style concurrent queue which is a common data structure encountered in the session type literature. A queue is described by the following type:

$$queue_A := \&\{ins : A \multimap queue_A, del : \oplus\{none : 1, some : A \otimes queue_A\}\}$$

The following diagram illustrates the channel topology of a client interacting with a queue server.



Each of the p_i nodes here represents a queue cell which holds a value and are linked together by bidirectional channels of type queue_A. As indicated by the type constructor &, the first queue node q_1 first receives either an ins or del label from the client. In the case of an ins label, p_1 receives a value v of type A (indicated by \multimap) from the client. The p_1 node then sends an ins label to p_2 and forwards v to it. This forwarding process repeats until the value reaches the end of the queue where a new queue cell p_{n+1} is allocated to store v. On the other hand, if p_1 receives a del label, the type constructor \oplus requires that p_1 send either none or some. The none label is sent to signify that the queue is empty and ready to terminate (indicated by 1). The some label is sent along with a value of type A (indicated by \otimes) which is the dequeued element. Finally, p_1 forwards its channel, connecting to p_2 , to the client so that the client may continue interacting with the rest of the queue.

It is clear from the example above that the session type $queue_A$ only lists what operations a queue should support, but does not specify the expected behavior of these operations. For instance, it does not specify that an ins operation should add an element to the back of the queue or that a

del operation should return the element at the front of the queue. A correct implementation needs to maintain all of these additional invariants not captured by the session type. In fact, due to the under specification of the queue $_A$ type, it is possible to implement a "queue" which simply ignores all ins messages and always returns none on del.

To address this issue, we develop TLL_C , a dependent session type system which extends the Two-Level Linear dependent type theory (TLL) [18] with session-typed concurrency. In TLL_C , one could define the queues through the following dependent session type:

```
queue(xs: list A) := ?(\ell: opr).match \ell with

| ins(v) \Rightarrow queue(snoc(xs, v))

| del \Rightarrow match xs with (x:: xs') \Rightarrow !(sing x).!(hc\(queue(xs')\(\rangle).1 | [] \Rightarrow 1
```

Here, the type queue(xs) is parameterized by a list xs which represents the current contents of the queue. Notice that the type no longer needs the \oplus and & type constructors to describe branching behavior. Instead, it uses type-level pattern matching to inspect the label ℓ received from the client. The opr type which ℓ inhabits is defined as a simple inductive type with two constructors:

```
inductive opr := ins : A \rightarrow \text{opr} \mid \text{del} : \text{opr}
```

When a queue server receives an ins(v) value, the type of the server becomes queue(snoc(xs, v)) were snoc appends v to the end of xs. Conversely, when a del label is received, the type-level pattern matching on xs enforces that if the queue is non-empty (i.e. x::xs' case), then the server must send the front element x of the queue to the client (indicated by the $singleton\ type\ sing\ x$) along with the channel $\mathbf{hc}(queue(xs'))$ connecting to the remainder of the queue. If the queue is empty (i.e. [] case), then the server simply terminates.

Given the queue protocol describe above, we can construct queue process nodes and interact with them. The following signatures are of helper functions that wrap interactions with the queue nodes into a convenient interface:

```
insert : \forall \{xs : \text{list } A\} \ (x : A) \rightarrow \text{Queue}(xs) \rightarrow \text{Queue}(\text{snoc}(xs, x))
delete : \forall \{x : A\} \ \{xs : \text{list } A\} \rightarrow \text{Queue}(x :: xs) \rightarrow C(\text{sing } x \otimes \text{Queue}(xs))
free : \text{Queue}([]) \rightarrow C(\text{unit})
```

The Queue type here is a type alias for the *channel type* of queues (explained later in detail) and the *C* type constructor here is the *concurrency monad* which encapsulates concurrent computations. Notice in the signature of insert and delete that there are dependent quantifiers surrounded by curly braces. These are the *implicit* quantifiers of TLL which indicate that the corresponding arguments are "ghost" values used for type checking and erased prior to runtime. For our purposes here, such ghost values are especially useful for *relationally* specifying the expected behaviors of queue interactions in terms of sequential list operations. For instance, the signature of insert states that the queue obtained after inserting *x* is related to the original queue by the list operation snoc. Similarly, the signature of delete states that deleting from a non-empty queue returns the front element *x*. Even though neither of these *xs* ghost values exist at runtime, they *statically* ensure that concurrent processes implementing these interfaces behave like actual queues, i.e., are first-in-first-out data structures. In a later section we will show how a generalized map-reduce algorithm can be implemented and verified using similar techniques.

Integrating session typed based concurrency into TLL is non-trivial due to the fact that TLL is a dependently typed functional language. While prior works [19, 45] have successfully combined *classical* session types with functional languages, its is well known that classical session types do not easily support recursive session types [20] (needed to express our queue type). The main

issue is that classical session types are defined in terms of a *dual* operator which does not easily commute with recursive type definitions. The addition of arbitrary type-level computations through dependent types further complicates this matter. On the other hand, *intuitionistic* session types [12] eschew the dual operator and define dual *interpretations* of session types based their *left* or *right* sequent rules. Because intuitionistic session types do not rely on a dual operator, they are able to support recursive session types without commutativity issues. However, intuitionistic session types are often formulated in the context of process calculi without a functional layer. To enjoy the benefits of intuitionistic session types in a functional setting, we develop a novel form of intuitionistic session types where we separate the notion of *protocols* from *channel types*. The queue(*xs*) type from before is, in actuality, a protocol whereas $\mathbf{hc}\langle \text{queue}(xs) \rangle$ is a channel type. In general, a channel type is formed by applying the $\mathbf{ch}\langle \cdot \rangle$ and $\mathbf{hc}\langle \cdot \rangle$ type constructors to protocols. These constructors provide dual interpretations to protocols, allowing dual channels of the same protocol to be connected together. For example, !*A.P* would be interpreted dually as follows:

```
ch\langle !A.P \rangle (send message of type A)
hc\langle !A.P \rangle (receive message of type A)
```

Such channel types can be naturally included into the contexts of functional type systems without needing to instrument the underlying language into a sequent calculus formulation. We believe our treatment of intuitionistic session types is not specific to TLL_C and is widely applicable for integrating intuitionistic session types with other functional languages.

In order to show that TLL_C ensures communication safety, we develop a process calculus based concurrency semantics. Process configurations in the calculus are collections of TLL_C programs interconnected by channels. At runtime, individual processes are evaluated using the program semantics of base TLL. When two processes at opposing ends (i.e. dually typed) of a channel are synchronized and ready to communicate, the process level semantics transmits their messages across the channel. We study the meta-theory of TLL_C and prove that it is indeed sound at both the level of terms and at the level of process configurations.

All lemmas and theorems reported in the this paper are formalized in Rocq [36]. All examples can be compiled into C programs using our prototype compiler where concurrent processes are implemented using POSIX threads. The compiler implements advanced language features such dependent pattern matching and functional in-place programming [27] for linear types. Proofs, source code, and examples are available in our git repository¹.

In summary, we make the following contributions:

- We extend the Two-Level Linear dependent type theory (TLL) with session type based concurrency, forming the language of TLL_C. TLL_C inherits the strengths of TLL such as Martin-Löf style linear dependent types and the ability to control program erasure.
- We develop a novel formulation of intuitionistic session types through a clear separation
 of protocols and channel types. We believe this formulation to be widely applicable for
 integrating session types into other functional languages.
- We study the meta-theoretical properties of TLL_C . We show that TLL_C , as a term calculus, possesses desirable properties such as confluence and subject reduction and, as a process calculus, guarantees communication safety.
- The entire calculus, with its meta-theorems, is formalized in Rocq.
- We implement a prototype compiler which compiles TLL_C into safe and efficient C code.

¹TODO

2 Overview of Dependent Session Types

Session types in TLL_C are *minimalistic* in design and yet surprisingly expressive due to the presence of dependent types. Through examples, we provide an overview of how dependent session types facilitate certified concurrent programming in TLL_C .

2.1 Message Specification

An obvious, but important, use of dependent session types is the precise specification of message properties communicated between parties. This is useful in practical network systems where the content of messages may depend on the value of a prior request. Consider the following protocol:

```
!(sz: nat). ?(msg: bytes). ?\{sizeOf(msg) = sz\}. 1
```

This example showcases the main primitives for constructing dependent protocols in TLL_C : the !(x:A).B and ?(x:A).B protocol actions. The syntax of these constructs take inspiration from binary session types [19, 45] and label dependent session types [38], however the meaning of these constructs in TLL_C is subtly different. In prior works, the ! marker indicates that the channel is to send and the ? marker indicates that the channel is to receive. In TLL_C , neither marker expresses sending or receiving per se, but rather an abstract action that needs to be interpreted through a channel type. Hence, the description of the messaging protocol above is stated to be informal. To assign a precise meaning to the protocol, we need to view it through the lenses of channel types:

```
ch\langle !(sz : nat). ?(msg : bytes). ?\{sizeOf(msg) = sz\}. 1\rangle

hc\langle !(sz : nat). ?(msg : bytes). ?\{sizeOf(msg) = sz\}. 1\rangle
```

Here, these two channel types are constructed using *dual* channel type constructors: $\mathbf{ch}\langle\cdot\rangle$ and $\mathbf{hc}\langle\cdot\rangle$. The $\mathbf{ch}\langle\cdot\rangle$ constructor interprets! as sending and? as receiving while the $\mathbf{hc}\langle\cdot\rangle$ constructor interprets! as receiving and? as sending. In other words, dual channel types interpret protocol actions in opposite ways. These constructors act similarly to the duality of left and right rules for intuitionistic session types [12]. Unlike intuitionistic session types which require the base type system to be based on sequent calculus, our channel types can be integrated into the type systems of functional languages so long as linear types are supported.

2.2 Dependent Ghost Secrets

Dependent ghost messages have interesting applications when it comes to message specification. Consider the following encoding of a idealized Shannon cipher protocol:

```
H(E, D) := \forall \{k : \mathcal{K}\} \ \{m : \mathcal{M}\} \to D(k, E(k, m)) =_{\mathcal{M}} m (correctness property) \mathcal{E}(E, D) := !\{k : \mathcal{K}\}. !\{m : \mathcal{M}\}. !(c : C). !\{H(E, D) \times (c =_{C} E(k, m))\}. 1
```

Given public encryption and decryption functions $E: \mathcal{K} \times \mathcal{M} \to C$ and $D: \mathcal{K} \times C \to \mathcal{M}$ respectively, the protocol $\mathcal{E}(E,D)$ begins by sending ghost messages: key k of type \mathcal{K} and message m of type \mathcal{M} . Next, the ciphertext c of type C, indicated by round parenthesis, is actually sent to

the client. Finally, the last ghost message sent is a proof object witnessing the correctness property of the protocol: c is obtained by encrypting m with key k. Observe that for the overall protocol, only ciphertext c will be sent at runtime while the other messages (secrets) are erased. The Shannon cipher protocol basically forces communicated messages to always be encrypted and prevents the accidental leakage of plaintext.

It is important to note that ghost messages and proof specifications, by themselves, are *not* sufficient to guaranteeing semantic security. An adversary can simply use a different programming language and circumvent the proof obligations imposed by TLL_C . However, these obligations are useful in ensuring that honest parties correctly follow *trusted* protocols to defend against attackers. For example, in the Shannon cipher protocol above, an honest party is required by the type system to send a ciphertext that is indeed encrypted from the (trusted) algorithm E.

Another, more concrete, example of using ghost messages to specify secrets is the Diffie-Hellman key exchange [17] protocol defined as follows:

The DH protocol is parameterized by publicly known integers p and g. Without loss of generality, we refer to the message sender for the first row of the protocol as Alice and the message sender for the second row as Bob. From Alice's perspective, she first sends her secret value a as a dependent ghost message to initialize her half of the protocol. Next, her public value A is sent as a real message to Bob along with a proof that A is correctly computed from values p, g and a (using modular exponentiation powm). At this point, Alice has finished sending messages and waits for message from Bob to complete the key exchange. She first "receives" Bob's secret b as a ghost message which initializes Bob's half of the protocol. Later, Bob' public value b is received as a real message along with a proof that b is correctly computed from b, b0 and b0. Notice that between Alice and Bob, the only the real messages b1 and b2 will be exchanged at runtime. The secret values b3 and b4 and the correctness proofs are all ghost message that are erased prior to runtime. Basically, the DH protocol forces communication between Alice and Bob to be encrypted and maintain secrecy at runtime.

```
\operatorname{def} \operatorname{Alice} (a p q : \operatorname{int}) (c : \operatorname{ch} \langle \operatorname{DH}(p, q) \rangle)
                                                                                                                def Bob (b p g : int) (c : hc\langle DH(p, g) \rangle)
: C(unit) :=
                                                                                                                 : C(unit) :=
    let c \Leftarrow \mathbf{send} \ c \{a\} in
                                                                                                                      let \langle \{a\}, c \rangle \leftarrow \mathbf{recv} \ c in
    let c \Leftarrow \mathbf{send} \ c \ (\mathsf{powm}(q, a, p)) in
                                                                                                                      let \langle A, c \rangle \leftarrow \mathbf{recv} \ c in
    let c \Leftarrow \mathbf{send} \ c \ \{ \mathbf{refl} \} in
                                                                                                                      let \langle \{pf\}, c\rangle \leftarrow \mathbf{recv} \ c in
    let \langle \{b\}, c \rangle \leftarrow \mathbf{recv} \ c in
                                                                                                                      let c \Leftarrow \text{send } c \{b\} in
                                                                                                                      let c \Leftarrow \mathbf{send} \ c \ (\mathsf{powm}(q, b, p)) in
    let \langle B, c \rangle \leftarrow \mathbf{recv} \ c in
    let \langle \{pf\}, c\rangle \leftarrow \mathbf{recv} \ c in
                                                                                                                      let c \Leftarrow \mathbf{send} \ c \ \{\mathsf{refl}\} in
    close(c)
                                                                                                                      wait(c)
```

The DH key exchange protocol can be implemented through two simple monadic programs Alice and Bob as shown above. The C type constructor here is the concurrency monad for integrating the *effect* of concurrent communication with the *pure* functional core of TLL_C . There are two kinds of send (and respectively recv) operations at play here. The first kind, indicated by send c {v} is for sending a ghost message v on channel c. After type checking, these ghost sends are compiled to no-ops to that they do not participate in runtime communication. The second kind, indicated by send c (v), is for sending a real message v on channel c. These real sends are compiled to actual messages in the generated code. Finally, the close and wait operations synchronize the termination of the protocol. Notice that the duality of channel types $\mathbf{ch}\langle \mathsf{DH}(p,g)\rangle$ and $\mathbf{hc}\langle \mathsf{DH}(p,g)\rangle$ ensure that

every send in Alice is matched by a corresponding receive in Bob and vice versa. Moreover, Alice and Bob are enforced by the type checker to correctly carry out the Diffie-Hellman key exchange.

3 Relational Verification via Dependent Session Types

Earlier in the introduction section, we showed a sketch of how dependent session types can be used for certified concurrent programming through the example of a concurrent queue. In this section, we provide a detailed account of how we can use dependent session types to construct a generic map-reduce system. Similarly to the queue example, we will verify the correctness of the map-reduce system by relating it to sequential operations on trees.

3.1 Construction of Map-Reduce

Map-reduce is a commonly used programming model for processing large data sets in parallel. Initially, map-reduce creates a tree of concurrently executing workers as illustrated in Figure 1. The client partitions the data into smaller chunks and sends them to the leaf workers of the tree. Next, each leaf worker applies a user-specified function f to each of its received data chunks and sends the results to its parent worker. When an internal worker receives results from its children, it combines the results using another user-specified binary function g. This procedure continues until the root worker computes the final result and sends it back to the client. Due to the fact that workers without data dependencies can operate concurrently, the overall system can achieve significantly better performance than sequential implementations of the same operations.

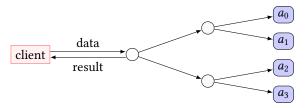


Fig. 1. Tree Diagram of Map-Reduce

The first step in constructing the map-reduce system is to build a model of our desired computation in a sequential setting. For this purpose, we define a simple binary tree inductive type:

```
inductive tree (A: \mathsf{U}) := \mathsf{Leaf} : A \to \mathsf{tree}(A) \mid \mathsf{Node} : \mathsf{tree}(A) \to \mathsf{tree}(A) \to \mathsf{tree}(A) def map : \forall \{A \ B : \mathsf{U}\} \ (f : A \to B) \to \mathsf{tree}(A) \to \mathsf{tree}(B) \mid \mathsf{Leaf} \ x \Rightarrow \mathsf{Leaf} \ (f \ x) \mid \mathsf{Node} \ l \ r \Rightarrow \mathsf{Node} \ (\mathsf{map} \ f \ l) (\mathsf{map} \ f \ r) def reduce : \forall \{A \ B : \mathsf{U}\} \ (f : A \to B) \ (g : B \to B \to B) \to \mathsf{tree}(A) \to B \mid \mathsf{Leaf} \ x \Rightarrow f \ x \mid \mathsf{Node} \ l \ r \Rightarrow g \ (\mathsf{reduce} \ f \ g \ l) \ (\mathsf{reduce} \ f \ g \ r)
```

In this definition, the type U of A is the universe of *unbound* (i.e. non-linear) types in TLL_C . So tree is parameterized by A which represents the type of data stored at the leaf nodes. The *sequential* map and reduce functions for tree are all defined in a standard way.

To construct the concurrent map-reduce system, the protocol of map-reduce must be able to branch depending on what operation the client requests to perform. Unlike many prior session type systems [12, 16] which provide built-in constructs (e.g. \oplus and &) for internal and external choice, we implement branching protocols using just dependent protocols and type-level pattern matching on sent or received messages. For our map-reduce system, we define the kinds of operations that can be performed through the inductive type opr:

```
\begin{array}{l} \operatorname{inductive} \operatorname{opr}(A:\mathsf{U}) := \operatorname{\mathsf{Map}} : \forall \{B:\mathsf{U}\} \ (f:A \to B) \to \operatorname{\mathsf{opr}}(A) \\ \mid \operatorname{\mathsf{Reduce}} : \forall \{B:\mathsf{U}\} \ (f:A \to B) \ (g:B \to B \to B) \to \operatorname{\mathsf{opr}}(A) \\ \mid \operatorname{\mathsf{Free}} : \operatorname{\mathsf{opr}}(A) \end{array}
```

The opr type has three constructors:

- Map f represents a map operation that applies the function f : A → B to each element of type A and produces results of type B.
- Reduce f g represents a reduce operation that first applies the function $f: A \to B$ to each element of type A and then combines the results using the binary function $g: B \to B \to B$.
- Free is the command that terminates the concurrent tree.

We are now ready to define the following treeP protocol to describe the interactions between nodes in the map-reduce tree.

```
def treeP (A : U) (t : \text{tree } A) := ?(o : \text{opr } A).

match o with Map _f f \Rightarrow \text{treeP } B \text{ (map } f t)

| \text{Reduce } _f g \Rightarrow !(\text{sing (reduce } f g t))}. \text{ treeP } t

| \text{Free} \Rightarrow \mathbf{1}
```

For each node n in the concurrent tree, it will be providing a channel of type $\operatorname{ch}\langle\operatorname{treeP} At\rangle$ to its parent. The parameter t of type tree A represents the shape of the sub-tree rooted at n. The treeP protocol states node n will receive a message o of type opr A from its parent. The protocol then branches, via type-level pattern matching on o, into three cases. If o is of the form Map f, then n will continue the protocol as treeP B (map f t). Notice that the type parameter of treeP is changed from A to B to reflect the fact that the data stored at the leaves of the sub-tree is transformed from type A to type B. Furthermore, the shape of the sub-tree has also changed from t to map t t. In the second case where t0 is of the form Reduce t1, t2, t3 will first send the result of type sing (reduce t3, t4) to its parent. The type sing t5 is the singleton type whose sole inhabitant is the element t5. After sending the result, t6 will continue the protocol as treeP t6, i.e. remains unchanged. Finally, t6 will terminate the protocol when t6 is Free.

Using the treeP protocol, we can now implement the worker processes that run at each node of the concurrent tree. The implementation of a leaf worker is shown below. We have elided uninteresting technical details regarding dependent pattern matching.

```
def leafWorker \{A: U\} (x:A) (c:\mathbf{ch}\langle \mathsf{treeP}\ A\ (\mathsf{Leaf}\ x)\rangle): C(\mathsf{unit}):= |\mathsf{et}\ \langle o,c\rangle := \mathbf{recv}\ c in match o with |\mathsf{Map}\Rightarrow \mathsf{leaf}\mathsf{Worker}\ \{B\}\ (f\ x)\ c |\mathsf{Reduce}\Rightarrow \mathsf{let}\ c\Leftarrow \mathbf{send}\ c\ (\mathsf{Just}\ (f\ x)) in leafWorker \{A\}\ x\ c |\mathsf{Free}\Rightarrow \mathbf{close}(c)
```

The leaf Worker function takes two non-ghost arguments: a data element x of type A and a channel c of type $\mathbf{ch}\langle \mathsf{treeP}\ A\ (\mathsf{Leaf}\ x)\rangle$. Through this channel c, the leaf worker will receive requests from its parent and provide responses accordingly. For instance, when the leaf worker receives a Map f request, it will apply $f:A\to B$ to its data element x and continue as a leaf worker with the new data element fx. In this case, the type parameter of leaf Worker has changed from A to B to reflect the transformation of the data element.

To represent internal node workers we implement the following nodeWorker function. This function takes (non-ghost) channels c_l and c_r of types $\mathbf{hc}\langle \text{treeP } A \ l \rangle$ and $\mathbf{hc}\langle \text{treeP } A \ r \rangle$ for communicating with its left and right children. Notice that the types of these channels are indexed by ghost values l and r of type tree A which represent the shapes of the concurrent sub-trees providing c_l

and c_r . The nodeWorker communicates with its parent through the channel c whose type is indexed by the ghost value Node l r.

```
def nodeWorker \{A : U\} \{l \ r : tree \ A\}
         (c_l : \mathbf{hc} \langle \mathsf{treeP} \ A \ l \rangle) \ (c_r : \mathbf{hc} \langle \mathsf{treeP} \ A \ r \rangle) \ (c : \mathbf{ch} \langle \mathsf{treeP} \ A \ (\mathsf{Node} \ l \ r) \rangle) : C(\mathsf{unit}) :=
    let \langle o, c \rangle := \mathbf{recv} \ c in
    match o with
      | Map f \Rightarrow
         let c_l \leftarrow \mathbf{send} \ c_l \ (\mathsf{Map} \ f) in
         let c_r \Leftarrow \mathbf{send} \ c_r \ (\mathsf{Map} \ f) in
         let c \Leftarrow \mathbf{send} \ c (Just unit) in
         nodeWorker \{B\} \{(\text{map } f \ l) \ (\text{map } f \ r)\}\ c_l\ c_r\ c
      | Reduce _f g \Rightarrow
         let c_l \Leftarrow \mathbf{send} \ c_l (Reduce f \ q) in
         let c_r \Leftarrow \mathbf{send} \ c_r \ (\mathsf{Reduce} \ f \ g) in
         let \langle \text{Just } v_l, c_l \rangle \leftarrow \mathbf{recv} \ c_l \text{ in}
         let \langle \text{Just } v_r, c_r \rangle \Leftarrow \text{recv } c_r \text{ in}
         let c \Leftarrow \mathbf{send} \ c \ (\mathsf{Just} \ (g \ v_l \ v_r)) in
         nodeWorker \{A\} \{l \ r\} c_l \ c_r \ c
      | Free \Rightarrow
         let c_l \Leftarrow \mathbf{send} \ c_l Free in
         let c_r \Leftarrow \mathbf{send} \ c_r Free in
         wait(c_l); wait(c_r); close(c)
```

Given the signature of nodeWorker and the definition of the treeP protocol, it is not hard to see that the implementation of nodeWorker is constrained to function exactly as intended. For instance, in the case where nodeWorker receives a Map f request from its parent, the type of c becomes $\mathbf{ch}\langle \text{treeP } B \text{ (map } f \text{ (Node } l \text{ } r))\rangle$ which simplifies to $\mathbf{ch}\langle \text{treeP } B \text{ (Node (map } f \text{ } l) \text{ (map } f \text{ } r))\rangle$. The shapes of the left and right sub-trees after the map operation need to become map f l and map l l respectively. In other words, the type of l forces the nodeWorker process to recursively send the Map l request to both of its children to transform them into sub-trees of type l l (map l l) and l l l (map l l).

3.2 A Certified Interface for Map-Reduce

Now that we have defined both leaf and internal node workers, we can wrap them up into a more convenient interface as presented below.

```
type cTree (A : U) (t : \text{tree } A) := C(\text{hc}(\text{treeP } t))

def cLeaf \{A : U\} (x : A) : \text{cTree } A (Leaf x) :=
	fork(c : \text{ch}(\text{treeP } A \text{ (Leaf } x))) with leaf Worker x c

def cNode \{A : U\}\{l \ r : \text{tree } A\} (c_l : \text{cTree } A \ l) (c_r : \text{cTree } A \ r) : \text{cTree } (\text{Node } l \ r) :=
	let c_l \Leftarrow c_l in
	let c_r \Leftarrow c_r in
	fork(c : \text{ch}(\text{treeP } A \text{ (Node } l \ r))) with nodeWorker c_l c_r c
```

The type alias cTree is defined to aid in the readability of the interface. The wrapper functions cLeaf and cNode respectively create leaf and internal node workers. This is accomplished by *forking* a new process using the **fork** construct of the concurrency monad. In particular, when given some a channel type $\mathbf{ch}\langle P\rangle$, the **fork** construct will create a new channel and give one end of it to the caller at type $\mathbf{hc}\langle P\rangle$ and spawn a new process that runs the worker with the other end of the channel at

type $\mathbf{ch}\langle P\rangle$. The duality of the channels types allows the caller and the worker to communicate. Using these wrapper functions, one can construct a concurrent tree in virtually the same way as one would construct a sequential tree. For example, the following code constructs a concurrent tree with four leaf nodes containing integers 0, 1, 2 and 3 respectively.

```
cNode (cNode (cLeaf 0) (cLeaf 1)) (cNode (cLeaf 2) (cLeaf 3))
```

The type of this expression is rather verbose to write manually as it contains the full shape of the concurrent tree. This is not a problem in practice as *constant* type arguments (such as the tree shapes here) can almost always be inferred automatically by the type checker.

Finally, we implement the cMap and cReduce functions that provide the map and reduce operations on concurrent trees. These functions are implemented by simply sending the appropriate requests to the root worker of the concurrent tree.

```
 \begin{aligned} & \operatorname{def} \operatorname{cMap} \left\{A \: B : \: \mathsf{U}\right\} \left\{t : \operatorname{tree} \: A\right\} \left(f : A \to B\right) \left(c : \operatorname{cTree} \: A \: t\right) : \operatorname{cTree} \: B \; (\operatorname{map} \: f \: t) := \\ & \operatorname{let} \: c \: \Leftarrow \: c \; \operatorname{in} \\ & \operatorname{let} \: c \: \Leftarrow \: \operatorname{send} \: c \; (\operatorname{Map} \: f) \; \operatorname{in} \\ & \operatorname{return} \: c \end{aligned}   \begin{aligned} & \operatorname{def} \: \operatorname{cReduce} \left\{A \: B : \: \mathsf{U}\right\} \left\{t : \operatorname{tree} \: A\right\} \left(f : A \to B\right) \left(g : B \to B \to B\right) \left(c : \operatorname{cTree} \: A \: t\right) := \\ & \operatorname{let} \: c \: \leftarrow \: \operatorname{c} \: \operatorname{in} \\ & \operatorname{let} \: c \: \Leftarrow \: \operatorname{send} \: c \; (\operatorname{Reduce} \: f \: g) \; \operatorname{in} \\ & \operatorname{let} \: \langle v, c \rangle \: \Leftarrow \: \operatorname{recv} \: c \operatorname{in} \\ & \operatorname{return} \; \langle v, \operatorname{return} \: c \rangle \end{aligned}
```

From the type signature of cMap, we can see that it takes a function f and a concurrent tree of type cTree A t and returns a new concurrent tree of type cTree B (map f t). In other words, the type of cMap guarantees that the shape of the concurrent tree is transformed in the same way as its sequential tree model under the map function. Similarly, the cReduce takes a concurrent tree of type cTree A t and returns a (linear) pair consisting of the result of type sing (reduce f g t), and the original concurrent tree. The correctness of cReduce is guaranteed by the singleton type of its result: reducing a concurrent tree results in the same value as reducing its sequential tree model.

3.3 Concurrent Mergesort via Map-Reduce

By properly instantiating the map-reduce interface defined previously, we can implement more complex concurrent algorithms. Moreover, dependent session types allows us to easily verify the correctness of these derived concurrent algorithms relationally through their sequential models. As an extended example, we implement a concurrent version of the mergesort algorithm using the map-reduce interface and verify its correctness.

We define sequential msort, as a model of our concurrent implementation, in the usual way using split and merge functions. We will not go into further details regarding the well-founded recursion of msort or the correctness of sorting as these are textbook results [14].

```
def split (xs : list int) : list int \times list int := ...

def merge (xs \ ys : list int) : list int := ...

def msort (xs : list int) : list int := match \ xs \ with

| \ nil \Rightarrow nil

| \ x :: nil \Rightarrow x :: nil

| \ zs \Rightarrow let \langle xs, ys \rangle := split \ zs \ in merge \ (msort \ xs) \ (msort \ ys)
```

Generally, to implement an algorithm using the map-reduce paradigm, one must first decompose the algorithm and data into a form that is amenable to parallelization. For mergesort, the input list can be recursively split into smaller sub-lists which can be processed in parallel. To make this decomposition *explicit*, we define the following splittingTree function that constructs a binary tree representation of how the input list is split by the mergesort algorithm.

```
def splittingTree (xs : list int) : tree (list int) := match xs with 
 | nil <math>\Rightarrow Leaf nil 
 | x :: nil \Rightarrow Leaf (x :: nil) 
 | zs \Rightarrow let \langle xs, ys \rangle := split zs in Node (splittingTree xs) (splittingTree ys)
```

To apply map-reduce, we need to construct a concurrent representation of its splitting tree with type cTree (list int) (splittingTree xs). While it is tempting to directly convert the result of splittingTree into a concurrent tree by recursively replacing Leaf with cLeaf and Node with cNode, such an approach would require traversing both the input list (to construct the splitting tree) and the resulting tree (to convert it into a concurrent tree). This would lead to a bottleneck in the performance of the overall algorithm as the traversals would be done sequentially without exploiting parallelism. Instead, we define the splittingCTree function that constructs the concurrent splitting tree in a concurrent manner.

```
def splittingCTree (xs : \text{list int}) : \text{ch} \langle !(\text{cTree (list int) (splittingTree } xs)). 1 \rangle \rightarrow C(\text{unit}) :=  match xs with | \text{nil} \Rightarrow \text{let } c \Leftarrow \text{send } c \text{ (cLeaf nil) in } \text{close}(c); \text{ return ()}  | x :: \text{nil} \Rightarrow \text{let } c \Leftarrow \text{send } c \text{ (cLeaf } (x :: \text{nil})) \text{ in } \text{close}(c); \text{ return ()}  | zs \Rightarrow \text{let } \langle xs, ys \rangle := \text{split } zs \text{ in}  | \text{let } c_l \Leftarrow \text{fork}(c) \text{ with } \text{splittingCTree } xs \text{ } c \text{ in}  | \text{let } c_r \Leftarrow \text{fork}(c) \text{ with } \text{splittingCTree } ys \text{ } c \text{ in}
```

The splittingCTree function takes an additional channel argument c which is used to send back the constructed concurrent tree to its caller. This small change allows the recursive case to fork two new processes to construct the left and right sub-trees in parallel. After both sub-trees have been constructed, the parent process can then combine them into a single concurrent tree using cNode and send it back to its caller. Notice that splittingCTree never calls the sequential splittingTree function and only uses it at the type level to model the concurrent tree being constructed. The complete implementation of splittingCTree can be found in the supplementary materials but is shortened here for brevity.

Now that we have constructed a concurrent splitting tree of our input list, we can apply the cReduce operation instantiated with $f := \lambda(x).x$ and g := merge to perform merging in parallel. This gives us an output of type

```
C(\text{sing (reduce }(\lambda(x).x) \text{ merge (splittingTree } xs))} \otimes \text{cTree (list int) (splittingTree } xs))
```

The singleton value sing (reduce $(\lambda(x).x)$ merge (splittingTree xs)) returned by the monad relationally describes this series of concurrent computations using just sequential operations. This allows us to easily verify the correctness of our concurrent mergesort implementation by proving the following theorem (in the internal logic of TLL) which states that reducing the splitting tree of a list is equivalent to performing mergesort on this list.

```
theorem reduceSplittingTree : \forall (xs : \text{list int}). \text{ reduce } (\lambda(x).x) \text{ merge (splittingTree } xs) = \text{msort } xs
```

Using this theorem, we can rewrite the singleton value returned by cReduce to sing (msort xs). In other words, the result of our concurrent mergesort implementation is guaranteed to be exactly the same as that of the sequential mergesort algorithm, thus completing our verification.

The full pipeline of concurrent mergesort is given in the following cMSort function.

```
def cMSort (xs: list int): C(\text{sing (msort } xs)) :=  let c \Leftarrow \text{fork}(c) with splittingCTree xs c in let \langle ctree, c \rangle \Leftarrow \text{recv } c in wait c; let \langle v, ctree \rangle \Leftarrow \text{cReduce } (\lambda(x).x) merge ctree in let ctree \Leftarrow \text{send } ctree Free in wait ctree; return (rewrite[reduceSplittingTree] v)
```

4 Formal Theory of Dependent Session Types

4.1 Core TLL

In this section, we give a brief summary of the Two-Level Linear dependent type theory (TLL) [18]. TLL is a dependent type theory that combines Martin-Löf-style dependent types [29] with linear types [21, 43]. Notably, TLL supports *essential linearity* [28] through the use of a stratified "two-level" typing system: the *logical* level and the *program* level. The typing judgments of the two levels are written and organized as follows:

provides types
$$\Gamma \vdash m : A \text{ (Logical Typing)} \qquad \qquad \Gamma; \Delta \vdash m : A \text{ (Program Typing)}$$
subjects to verify

First, the *logical* level is a standard dependent type system that supports unrestricted usage of types and terms. The primary purpose of the logical level is to provide typing rules for types which will be used at the logical level. For example, the rules for dependent function type (Π -types) formation are defined at the logical level as follows:

$$\frac{\Gamma \vdash A : s \qquad \Gamma, x : A \vdash B : r}{\Gamma \vdash \Pi_t(x : A).B : t} \qquad \frac{\prod_{P \vdash A : s} \Gamma, x : A \vdash B : r}{\Gamma \vdash \Pi_t\{x : A\}.B : t}$$

The symbols s, r, t range over the *sorts* of type universes, i.e. U or L. These sorts are used to classify types into two categories: unrestricted types (A : U) and linear types (A : L). Program level terms which inhabit unrestricted types can be freely duplicated or discarded, while those which inhabit linear types must be used exactly once. Note that this usage restriction is *not* enforced at the logical level as the logical level typing judgment is completely structural. This is safe because the logical level will never be executed at runtime and is only used for type checking and verification. Thus, multiple uses of a linear resource at the logical level will not lead to any runtime errors.

At the program level, the typing judgment Γ ; $\Delta \vdash m : A$ is used to exclusively type *terms*. In other words, no rules for forming types are defined at the program level. All the types used in Γ , Δ , m and A must be well-formed according to the logical level typing judgment. This typing judgment possesses two contexts: Γ of all variables in scope, and Δ of all variables that are computationally relevant in program m. Context Δ is crucial for enforcing linearity at the program level. For example, consider the λ -abstraction rules:

$$\frac{\Gamma, x: A; \Delta, x:_{s} A \vdash m: B}{\Gamma; \Delta \vdash \lambda_{t}(x:A).m: \Pi_{t}(x:A).B} \xrightarrow{\text{IMPLICIT-LAM}} \frac{\Gamma, x: A; \Delta \vdash m: B}{\Gamma; \Delta \vdash \lambda_{t}\{x:A\}.m: \Pi_{t}\{x:A\}.B}$$

In Explicit-Lam, we can see that the bound variable x is added to both contexts Γ and Δ . This indicates that x is a variable which can be used both logically (in types and ghost values) through Γ , and computationally (in real values) through Δ . On the other hand, in the Implicit-Lam rule, x is only added to Γ but not Δ . This indicates that x is a ghost variable which can only be used logically. A ubiquitous example of ghost variables are type parameters in polymorphic functions. For example, the polymorphic identity function can be implemented as

$$\lambda_{\mathsf{U}}\{A:\mathsf{U}\}.\lambda_{\mathsf{U}}(x:A).x$$

which has the type $\Pi_U\{A:U\}$. $\Pi_U(x:A)$.A. Arguments to implicit functions are typed at the logical level, thus allowing polymorphic functions to be instantiated with a type as an argument. Additionally, as demonstrated in the examples of prior sections, ghost variables also facilitate program verification by statically describing abstractions and invariants of program states.

In the two λ -abstraction rules above, the premise $\Delta \triangleright t$ is a simple side condition that states: if t=U, then all variables in Δ must be unrestricted. In other words, the λ -abstractions that can be applied unrestrictedly (with t=U) are not allowed to capture linearly typed variables from Δ . This is similar to the restriction imposed on closures implementing the Fn trait (i.e. those that can be called multiple times) in Rust [37] where capturing of mutable references is prohibited. If such a restriction is not imposed, then evaluating a λ -abstraction (that captures a linear variable) twice may lead to unsafe memory accesses such as double frees or use-after-frees.

The application rules for both explicit and implicit functions are as follows:

$$\frac{\Gamma; \Delta_1 \vdash m : \Pi_t(x : A).B \qquad \Gamma; \Delta_2 \vdash n : A}{\Gamma; \Delta_1 \cup \Delta_2 \vdash m \; n : B[n/x]} \frac{\Gamma; \Delta_1 \cup \Delta_2 \vdash m \; n : B[n/x]}{\Gamma; \Delta \vdash m \; \{n\} : B[n/x]}$$

$$\frac{\Gamma; \Delta_1 \cup \Delta_2 \vdash m \; n : B[n/x]}{\Gamma; \Delta \vdash m \; \{n\} : B[n/x]}$$

In Explicit-App, the argument n is a real value which must be typed at the program level. The \cup operator merges the two program context Δ_1 and Δ_2 by contracting unrestricted variables and requiring that linear variables be disjoint, thus preventing the sharing of linear resources. In Implicit-App, the argument n is a ghost value that is typed at the logical level. Due to the fact that ghost values are erased prior to runtime, the program context Δ in the conclusion only tracks the computationally relevant variables used in m. Notice how in Explicit-App, the argument n is substituted into the return type B. This allows types to depend on program level terms regardless of whether they are of linear or unrestricted types.

Usage vs Uniqueness. Compared to other linear dependent type theories [5, 13, 28, 30, 42] which only enforce the linear *usage* of resources, the TLL type system prevents the *sharing* of linear resources as well. This is similar to the subtle distinction between linear logic [21] and bunched implications [31, 32] described by O'Hearn. Consider a linear function f, in the aforementioned dependent type theories, of some type $A \multimap B$. When function f is applied to some argument v of type A, the argument v is guaranteed to be used exactly once in the *body* of f. Notice that this notion of linearity does not guarantee that f has unique access to v. If v was obtain from some !-exponential or ω-quantity (the sharable quantity in graded systems [5, 30]), then there may be other aliases of v which can be used outside of f.

Wadler, in his seminal work [44], made a similar distinction between linearity and uniqueness in the context of functional programming, noting that implicit uses of *promotion* and *dereliction* in linear logic can lead to violations of uniqueness. He coins the term *steadfast types* to refer to type systems that enforce both linearity and uniqueness. In this sense, TLL is steadfast as its *sort-uniqueness* property (i.e. types uniquely inhabit either U or L) prohibits the implicit promotion and dereliction of linear types, thus preventing the sharing of linear resources. The heap semantics [41] of TLL shows that its programs enjoy the *single-pointer* property which is a consequence of uniqueness at

runtime. In the context of concurrency, the steadfast type system of TLL makes it especially suitable for integration with session types: linear usage prevents replaying of communication protocols and uniqueness ensures that a communication channel has a single owner.

4.2 Dependent Session Types of TLL_C

In this section, we formally present the dependent session types of TLL_C .

Basic Protocols and Channel Types. The intuitionistic session types of TLL_C are decoupled into *protocols* and *channel types*. The rule for forming protocols is as follows:

Ргото	Explicit-Action	IMPLICIT-ACTION	End
Γ ⊢	$\Gamma, x : A \vdash B : \mathbf{proto}$	$\Gamma, x : A \vdash B : \mathbf{proto}$	Γ \vdash
$\Gamma \vdash \mathbf{proto} : U$	$\Gamma \vdash \rho(x:A). B : \mathbf{proto}$	$\Gamma \vdash \rho\{x:A\}.\ B: \mathbf{proto}$	$\overline{\Gamma \vdash 1 : \mathbf{proto}}$

where
$$\rho \in \{!,?\}$$

Here, the Proto rule introduces the **proto** type which is the type of all protocols. Note that **proto** is an unrestricted type, thus protocols can be freely duplicated or discarded. The Explicit-Action and Implicit-Action rules form dependent protocols which inhabit the **proto** type. The End rule marks the termination of a protocol.

Once a protocol is defined, we can form channel types using the following rules:

СнТүре	НсТүре	
$\Gamma \vdash A : \mathbf{proto}$	$\Gamma \vdash A : \mathbf{proto}$	
$\Gamma \vdash \mathbf{ch}\langle A \rangle : L$	$\Gamma \vdash \mathbf{hc}\langle A \rangle : L$	

Notice that the channel type constructors $\mathbf{ch}\langle\cdot\rangle$ and $\mathbf{hc}\langle\cdot\rangle$ lift protocols, which are unrestricted values, into linear types. This means that channels must be used exactly once. Furthermore, as explained in the previous section, the unique ownership of linear types in TLL ensures that only a single entity has access to a channel at any point in time, thus preventing race conditions.

Recursive Protocols. Recursive protocols can be formed using the $\mu(x:A)$.m construct:

FIXPOINT $\Gamma, x: A \vdash m: A$ A is an arity ending on **proto** x is guarded by protocol action in m $\Gamma \vdash \mu(x:A).m:A$

For a $\mu(x:A).m$ term, we require that A be an arity ending on **proto**. This prevents μ from introducing logical inconsistencies as it can only be used to construct protocols and not proofs for arbitrary propositions. To ensure that protocols defined through $\mu(x:A).m$ can be productively unfolded, recursive usages of x must be syntactically guarded behind a protocol action in m. This enforces the contractiveness condition for recursive session types [19]. Both the arity and guardedness conditions are stable under substitution. Due to space limitations, we present the rules of arities and guardedness in the appendix.

The difficulty of integrating recursive protocols in classical session type systems is well documented [20]. The key challenge is to define a suitable *duality* operator that commutes with recursion. The following example is due to Bernardi and Hennessy [8]. Suppose we define a reasonable, but naive, duality operator $(\cdot)^{\perp}$ which simply flips! and? in protocols. For the dual of recursive protocol $\mu X.?X.X$, if we first apply duality and then unfold the recursion, we get:

$$(\mu X.?X.X)^{\perp} = \mu X.!X.X = !(\mu X.!X.X).(\mu X.!X.X)$$

On the other hand, if we first unfold the recursion and then apply duality, we get:

$$(\mu X.?X.X)^{\perp} = (?(\mu X.?X.X).(\mu X.?X.X))^{\perp} = !(\mu X.?X.X).(\mu X.!X.X)$$

Notice that the resulting protocols do not agree on the type of the sent message. While solutions have been proposed to address this issue [8, 9], they do not generalize to dependent session types due to the presence of arbitrary type-level computation. In TLL_C , the separation of protocols and channels types allows us to sidestep the duality problem entirely. Suppose we define our previously problematic recursive protocol in TLL_C as follows:

$$T \triangleq \mu(X : proto).?(_:X).X = ?(_:\mu(X : proto).?(_:X).X). \ \mu(X : proto).?(_:X).X$$

When viewed through the lens of channel type constructors $\mathbf{ch}\langle\cdot\rangle$ and $\mathbf{hc}\langle\cdot\rangle$, the actions specified by the unfolded protocol are correctly dual to each other. More specifically, a channel of type $\mathbf{ch}\langle T\rangle$ receives a protocol of type T whereas a channel of type $\mathbf{hc}\langle T\rangle$ sends a protocol of type T.

Concurrency Monad. Concurrency is integrated into the pure functional core of TLL through a concurrency monad *C*. The basic components of the monad are given in the following rules.

$$\begin{array}{ccc} C_{\text{TYPE}} & & & & & & & & \\ \Gamma \vdash A : s & & & & & & \\ \hline \Gamma \vdash C(A) : \mathsf{L} & & & & & \\ \hline \Theta; \Gamma; \Delta \vdash \mathbf{return} \, m : C(A) & & & & & \\ \hline \Theta; \Gamma; \Delta \vdash \mathbf{return} \, m : C(A) & & & & & \\ \hline \Theta_1 \cup \Theta_2; \Gamma; \Delta_1 \cup \Delta_2 \vdash \mathbf{let} \, x \Leftarrow m \, \mathbf{in} \, n : C(B) & & \\ \hline \Theta_1 \cup \Theta_2; \Gamma; \Delta_1 \cup \Delta_2 \vdash \mathbf{let} \, x \Leftarrow m \, \mathbf{in} \, n : C(B) & & \\ \hline \end{array}$$

To reason about the communication channels that will appear at *runtime*, the program level typing judgment is extended to include a *channel context* Θ which tracks the channels used by the program. It is crucial to understand that the channel context is largely a technical device for analyzing the type safety of TLL_C . Prior to runtime, the channel context is empty as no channels have been created. Programming is carried out using normal variables in Δ . At runtime, channels will be created and substituted for appropriate variables in Δ . It is these runtime channels that occupy the channel context Θ and are typed as follows:

$$\begin{array}{c} \text{Channel-CH} \\ \underline{\Gamma; \Delta \vdash \quad \epsilon \vdash A : \textbf{proto} \quad \Delta \triangleright \textbf{U}} \\ \hline c :_{\textbf{L}} \ \textbf{ch} \langle A \rangle; \Gamma; \Delta \vdash c : \textbf{ch} \langle A \rangle \\ \end{array}$$

The protocol A used in the channels types here must be *closed*. This is because channels at runtime must follow fully concretized protocols. The Γ and Δ contexts are allowed to be non-empty for the purely technical reason of facilitating proofs for renaming and substitution lemmas.

As explained in Section 2.1, the protocol actions !(x : A).B and ?(x : A).B are abstract constructs that need to be interpreted through channel types. Since $\mathbf{ch}\langle\cdot\rangle$ and $\mathbf{hc}\langle\cdot\rangle$ interpret protocol actions in opposite ways, we only present the typing rules for $\mathbf{ch}\langle\cdot\rangle$ below.

$$\begin{array}{ll} \text{Explicit-Send-CH} & \text{Explicit-Recv-CH} \\ \Theta; \Gamma; \Delta \vdash m : \textbf{ch} \langle !(x : A) . B \rangle & \Theta; \Gamma; \Delta \vdash m : \textbf{ch} \langle ?(x : A) . B \rangle \\ \hline \Theta; \Gamma; \Delta \vdash \textbf{send} \ m : \Pi_{\mathsf{L}}(x : A) . C(\textbf{ch} \langle B \rangle) & \Theta; \Gamma; \Delta \vdash \textbf{recv} \ m : C(\Sigma_{\mathsf{L}}(x : A) . \textbf{ch} \langle B \rangle) \\ \hline \text{Implicit-Send-CH} & \text{Implicit-Recv-CH} \\ \Theta; \Gamma; \Delta \vdash m : \textbf{ch} \langle !\{x : A\} . B \rangle & \Theta; \Gamma; \Delta \vdash m : \textbf{ch} \langle ?\{x : A\} . B \rangle \\ \hline \Theta; \Gamma; \Delta \vdash \textbf{send} \ m : \Pi_{\mathsf{L}}\{x : A\} . C(\textbf{ch} \langle B \rangle) & \Theta; \Gamma; \Delta \vdash \textbf{recv} \ m : C(\Sigma_{\mathsf{L}}\{x : A\} . \textbf{ch} \langle B \rangle) \\ \hline \end{array}$$

For the Explicit-Send-CH rule, a channel of type $\mathbf{ch}\langle !(x:A).B\rangle$ is applied to the **send** operator. This produces a function which takes a real value v of type A and returns a concurrent computation of type $C(\mathbf{ch}\langle B[v/x]\rangle)$ which represents the continuation of the protocol after sending a real value

of type A. When this monadic value is bound by rule BIND and executed at runtime, the value v will be sent on channel m. The dual Explicit-Recv-HC rule, as shown here,

Explicit-Recv-HC
$$\frac{\Theta; \Gamma; \Delta \vdash m : \mathbf{hc} \langle !(x : A). B \rangle}{\Theta; \Gamma; \Delta \vdash \mathbf{recv} \ m : C(\Sigma_{\mathsf{L}}(x : A). \mathbf{hc} \langle B \rangle)}$$

receives on a channel of type $\mathbf{hc}\langle !(x:A).B\rangle$ which produces a (monadic) dependent pair (similarly to Explicit-Recv-CH). The first component of the pair is the value of type A that was received, and the second component is a channel of type $\mathbf{hc}\langle B[v/x]\rangle$ representing the continuation of the protocol. Notice that, due to the linearity of the C monad, all of the intermediate monadic values are guaranteed to be bound by the BIND rule and executed.

The implicit send and receive rules are similar to their explicit counterparts, except that they send and receive ghost values instead of real values. This distinction manifests by having the <u>send</u> and <u>recv</u> operators produce implicit functions and implicit pairs respectively. When the implicit function of Implicit-Send-CH is applied to a ghost argument using Implicit-App (Section 4.1), the ghost argument will be erased prior to runtime. Similarly, the first component of the implicit pair produced by Implicit-Recv-CH is also an erased ghost value. The underlying type system of TLL ensures that these ghost values will only be used logically, thus are safe to erase.

The last communication rules govern the creation and termination of channels:

CLOSE and WAIT are simple rules used to free channels whose protocols have terminated. The FORK rule is used for creating a child process which concurrently executes the monadic computation m. The child process is provided with a fresh channel of type $\mathbf{ch}\langle A\rangle$ which is bound to the variable x in m. Dually, the parent process obtains the channel endpoint of type $\mathbf{hc}\langle A\rangle$, which can be used to communicate with the spawned process. Note that the newly spawned process m is allowed to capture pre-existing channels from Θ and program variables from Δ . Compared to intuitionistic session type systems based on the sequent calculus [12, 16, 33], the $\mathbf{ch}\langle A\rangle$ channel handed to the child process behaves like the right-hand side of a sequent (i.e. the *provided* channel), while the $\mathbf{hc}\langle A\rangle$ channel handed to the parent process behaves like the left-hand side of a sequent (i.e. the *consumed* channels). Essentially, we have embedded intuitionistic session types into a functional language without needing to reorganize the underlying type system into a sequent calculus formulation.

5 Semantics and Meta-Theory

5.1 Process Configurations

In the previous section, we have presented the typing rules for TLL_C terms which form individual processes. To compose multiple processes together, we introduce the process level typing judgment $\Theta \Vdash P$ below. This judgment formally states that a configuration of processes P is well-typed under the context Θ , which tracks the channels used by the processes in P at runtime.

$$\frac{\text{Expr}}{\Theta; \epsilon; \epsilon \vdash m : C(\mathsf{unit})} \qquad \frac{\text{Par}}{\Theta_1 \Vdash P_1} \qquad \frac{\text{Scope}}{\Theta_2 \Vdash P_2} \qquad \frac{\Theta, c :_{\mathsf{L}} \mathsf{ch}\langle A \rangle, d :_{\mathsf{L}} \mathsf{hc}\langle A \rangle \Vdash P}{\Theta \Vdash vcd.P}$$

The process configuration rules are standard. The Expr rule lifts well-typed closed terms of type C(unit) to processes. It is important for the term m to be closed as processes in a configuration cannot rely on external substitutions to resolve free variables, they can only communicate through

channels. In the PAR rule, well-typed configurations P and Q can be composed in parallel as long as their contexts Θ_1 and Θ_2 can be combined. The Scope rule allows two dual channels to be connected together, allowing processes holding channels c and d to communicate.

The structural congruence of process configurations is defined as the least congruence relation generated by the following standard rules:

$$P \mid Q \equiv Q \mid P$$
 $O \mid (P \mid Q) \equiv (O \mid P) \mid Q$ $P \mid \langle \mathbf{return} () \rangle \equiv P$ $vcd.P \mid Q \equiv vcd.(P \mid Q)$ $vcd.P \equiv vdc.P$ $vcd.vc'd'.P \equiv vc'd'.vcd.P$

Structural congruence states that parallel composition is commutative and associative and compatible with channel scoping. Processes which terminate with the unit value () can be removed from a configuration. Intuitively, two structurally congruent configurations should be considered equivalent regarding their communication behavior.

5.2 Semantics

Term Reduction. The operational semantics of TLL_C programs is mostly the same as that of call-by-value TLL [18]. The relation $m \rightsquigarrow m'$ is used to denote a single step of *program* level reduction. Due to the monadic formulation of concurrency in TLL_C , the only additional (non-trivial) program reduction rule is the following BINDELIM rule which reduces a monadic **let**-expression when its bound term is a **return** expression:

(BINDELIM) **let**
$$x \leftarrow \text{return } v \text{ in } m \rightsquigarrow m[v/x]$$
 (where v is a value)

Values now additionally include channels, partially applied communication operators and thunked monadic expressions. We will use the metavariable v to denote values for the rest of this paper. The full definition of values is presented in the appendix.

Process Reduction. The semantics of processes is defined through the relation $P \Rightarrow Q$ which states that process configuration P reduces to process configuration Q in one step. The process reduction rules are presented below.

The first four rules define the synchronous communication semantics of TLL_C . The Proc-Fork rule creates a pair of dual channels c and d to connect the continuation n of the parent process with the newly forked child process m. We can see here that the newly created channels c and d are substituted for the variables x and y in n and m respectively.

The Proc-End rule synchronizes the termination of communicating on dual channels c and d. The resulting process configuration contains two processes which are no longer connected by any channels. Additionally, the close and wait operations are replaced by unit return values once the termination is synchronized.

The Proc-Com rule governs the communication of a real message v from a sender to a receiver. The sending process continues as m with the channel c while the receiving process continues as n with the received message v and the channel d paired together as $\langle v, d \rangle_{L}$.

The Proc-Com rule is similar to Proc-Com except that it handles the communication of a ghost message o. While this rule seems to indicate that ghost messages are communicated at runtime, we will later show through the erasure safety theorem that ghost messages are always safe to be erased. The exchange of ghost messages here is only for the purpose of establishing a reference point for reasoning about the correctness of erasure safety.

The remaining four rules are standard. The Proc-Expr rule allows a singleton process to reduce by reducing its underlying term. The Proc-Par and Proc-Scope rules allow a process to reduce in parallel composition and under channel scope respectively. Finally, the Proc-Congr rule allows processes to reduce up to structural congruence.

5.3 Meta-Theory

Compatibility. We first show that the concurrency extensions of TLL_C are compatible with the underlying TLL type system. To this end, we prove that TLL_C enjoys the same meta-theoretical properties as TLL. Due to the fact that these properties do not involve concurrency, their proofs indicate that TLL_C is sound as a term calculus. Here we present a few representative theorems. The full list of theorems and their proofs can be found in our Rocq formalization.

The first theorem we present is the validity theorem which states that well-typed terms have well-sorted types. This theorem is important as it ensures that the types appearing in typing judgments are indeed valid (i.e. they inhabit a sort).

```
THEOREM 5.1 (VALIDITY). Given \Theta; \Gamma; \Delta \vdash m : A, there exists sort s such that \Gamma \vdash A : s.
```

In TLL and TLL_C , the sort of a type determines whether the type is a unrestricted or linear. This means that it is crucial for a type to have a unique sort, otherwise the same type could be interpreted as both unrestricted and linear, leading to unsoundness. To address this concern, we prove the sort uniqueness theorem below which states that a type can have at most one sort. This ensures no ambiguity on whether a type is to be considered unrestricted or linear.

```
THEOREM 5.2 (SORT UNIQUENESS). Given \Gamma \vdash A : s and \Gamma \vdash A : t, we have s = t.
```

The next theorem we present is the standard subject reduction theorem which states that types are preserved under term reduction. This theorem is necessary for ensuring that session fidelity holds during process reduction as singleton processes reduce by reducing their underlying terms.

```
THEOREM 5.3 (SUBJECT REDUCTION). Given \Theta; \epsilon; \epsilon \vdash m : A and m \rightsquigarrow m', we have \Theta; \epsilon; \epsilon \vdash m' : A.
```

Session Fidelity. The session fidelity theorem ensures that processes adhere to the communication protocols specified by their types. This property guarantees that well-typed processes will not encounter communication mismatches at runtime. Since we consider processes up to structural congruence, we must first show that configuration typing is preserved under structural congruence. This manifests as the following lemma.

```
Lemma 5.4 (Congruence). Given \Theta \Vdash P and P \equiv Q, we have \Theta \Vdash Q.
```

The session fidelity theorem is then stated as follows.

```
Theorem 5.5 (Session Fidelity). Given \Theta \Vdash P and P \Rightarrow Q, we have \Theta \Vdash Q.
```

One of the primary challenges in proving session fidelity is to show that typing is preserved during communication steps, specifically the Proc-Com, and Proc- $\underline{\text{Com}}$ cases. In these cases, the

message being communicated is transported from the sender to the receiver without the use of a substitution. We need to show that the message, after communication, is consistently typed with regards to the receiver's context. Unlike simple type systems where one could simply place a value into any context so long as the value has the expected type, dependent type systems require more care. For instance, the evaluation context $\langle [\cdot], \text{refl} \rangle : \Sigma(x : \text{nat}).(x = 1)$ is well-typed if and only if the hole is filled with 1. To address this challenge, we design the monadic BIND rule (Section 4.2) to disallow dependency on the bound value. More specifically, for **let** $x \leftarrow m$ **in** n expressions, the type of n cannot depend on x. This restriction means that m can be replaced by any other expression of the same type without affecting the type of n. Consider the Proc-Com step below:

```
vcd.(\langle \text{let } x \Leftarrow \text{send } c \text{ } v \text{ in } m \rangle \mid \langle \text{let } y \Leftarrow \text{recv } d \text{ in } n \rangle)

\Rightarrow vcd.(\langle \text{let } x \Leftarrow \text{return } c \text{ in } m \rangle \mid \langle \text{let } y \Leftarrow \text{return } \langle v, d \rangle_{\mathsf{L}} \text{ in } n \rangle)
```

This operation is carried out between two singleton processes that are evaluating monadic **let**-expressions. Due to the dependency restriction of the BIND rule, we can replace **send** c v with **return** c and **recv** d with **return** $\langle v, d \rangle_{L}$ without affecting the types of m and n. Due to the fact that all communication operations in TLL_{C} are carried out on **let**-expressions, the dependency restriction ensures that session fidelity holds during communication steps.

Global Progress. Global progress, i.e. deadlock-freedom, is a desirable property for concurrent programs. Many session type systems [12, 16, 45] guarantee global progress by construction through a disciplined use of channels. However, there are also session type systems [6, 23, 24, 38] that eschew global progress in favor of more expressive session types. TLL_C belongs to the latter category if we consider arbitrary well-typed process configurations. This is because the process type system of TLL_C does not prevent cyclic channel topologies that can lead to deadlocks. However, we can still prove a weaker form of global progress for TLL_C by considering only *reachable* process configurations. Intuitively, reachable configurations are those that can be constructed by **fork** operations starting from a single process. The global progress theorem is then stated as follows.

Theorem 5.6 (Global Progress). Given $\epsilon \Vdash P$ where P is reachable, either

- $P \equiv \langle return() \rangle$, or
- there exists Q such that $P \Rightarrow Q$.

Erasure Safety. To show that ghost messages are safe to erase, we define an erasure relation Θ ; Γ ; $\Delta \vdash m \sim m' : A$. This relation states that all ghost arguments and type annotations in m are replaced by a special opaque value \Box in m'. This relation is similar to the one defined for the erasure of *propositions* in standard dependent type theories [7, 26, 35]. The most important erasure rule is shown below. The full set of erasure rules can be found in the appendix.

$$\frac{\Theta; \Gamma; \Delta \vdash m \sim m' : \Pi_t\{x : A\}.B \qquad \Gamma \vdash n : A}{\Theta; \Gamma; \Delta \vdash m \{n\} \sim m' \{\square\} : B[n/x]}$$

The Erase-Implicit-App rule states that when erasing an implicit application m $\{n\}$, the ghost argument n is replaced by \square in the erased term. Consider the <u>send</u> c operator for sending ghost messages on channel c. As defined in Section 4.2, this partially applied operator has a type of the form $\Pi_L\{x:A\}.C(B)$. When fully applied as <u>send</u> c $\{n\}$, the ghost argument n is erased to \square by Erase-Implicit-App. Since \square is an opaque value, it cannot be inspected or pattern matched on. Thus, if programs can be evaluated soundly after erasing all ghost arguments and type annotations, we can conclude that ghost messages are safe to erase.

The erasure relation is then naturally lifted to the process level as $\Theta \Vdash P \sim P'$ where P' is the erased version of P. The rules for this relation are as follows:

We show that erasure is safe through the following two theorems. These theorems tell us that any possible reduction on an original object (either a term or process) can be simulated on its erased counterpart. Moreover, the erased object obtained after reduction also satisfies the erasure relation with respect to the reduced original object. Basically, these theorems state that any possible evaluation path of the original object remains valid after erasure.

THEOREM 5.7 (TERM SIMULATION). Given Θ ; ϵ ; $\epsilon \vdash m \sim m' : A$ and $m \rightsquigarrow n$, there exists n' such that $m' \rightsquigarrow^* n'$ and Θ ; ϵ ; $\epsilon \vdash n \sim n' : A$.

Theorem 5.8 (Process Simulation). Given $\Theta \Vdash P \sim P'$ and reduction $P \Rightarrow Q$, there exists Q' such that $P' \Rightarrow Q'$ and $\Theta \Vdash Q \sim Q'$.

6 Implementation

We implement a prototype compiler for TLL_C . The main components of the compiler are written in OCaml while a minimalistic runtime library is implemented in C. The compiler takes TLL_C source files as input and generates safe C code which can be further compiled into executable binaries on POSIX compliant systems. In this section, we give an overview of the inference, linearity checking and optimization phases of the compiler.

Inference. To reduce code duplication and type annotation burden, we implement two forms of inference: (1) automatic instantiation of *sort-polymorphic schemes* similarly to the TLL compiler and (2) elaboration of inferred arguments. Consider the identity function below:

$$\mathsf{def} \; \mathsf{id} \langle s \rangle \; \% \{A : \mathsf{Type} \langle s \rangle \} \; (x : A) : A := x$$

This function is a sort-polymorphic scheme as it is parameterized over sort variable s. Depending on the universe of A, sort s can be instantiated to either L for linear types or U for unrestricted types. This eliminates the need to define two separate identity functions for linear and unrestricted types. The type A here is marked by % to indicate that it is an inferred argument. Suppose id is applied to a natural number 42, the compiler creates two metavariables \hat{s} and $\hat{\alpha}$ to represent the elided sort and type arguments respectively. Type inference produces the following constraints:

$$id \ 42 \xrightarrow{\text{desugar}} id \ \ \hat{s} > \{\hat{\alpha}\} \ 42 \xrightarrow{\text{infer}} \begin{cases} \hat{s} = \mathsf{U} & \xrightarrow{\text{mono}} \\ \hat{\alpha} = \mathsf{Nat} \end{cases} id \ \ \mathsf{V} > \{\mathsf{Nat}\} \ 42$$

Once the constraints are solved through unification [1], the metavariables are replaced by their solutions. The monomorphized code is then passed to the next phase for linearity checking.

Linearity Checking. During the inference phase, the usage of linear variables is not tracked. The type checking algorithm essentially treats TLL_C as a fully structural type system. It is only after all sort-polymorphic schemes and inferred arguments are instantiated that the linearity checking begins. A substructural type checking algorithm is applied to determine if the elaborated program compiles with the actual typing rules of TLL_C . We adopt this two-phase approach to simplify the linearity checking algorithm. Although sort-polymorphism greatly reduces code duplication from the user's perspective, it also obfuscates the classification of types into linear and unrestricted ones. Thus, it is much easier to check linearity after monomorphization.

To support dependent pattern matching, we implement a variation of Cockx's algorithm [15] to type check **match**-expressions and elaborate them into well-formed case trees. Cockx's algorithm

forms the basis of Agda's [4] pattern matching mechanism. Several extensions are made to the original algorithm to account for pattern matching on linear inductive types and ghost terms. Our modified algorithm is able to correctly track resource usage in subtle cases such as nested patterns involving linear inductive types. We plan to present the details of our algorithm in a future publication.

Optimization. Once linearity checking is complete, ghost terms are erased in a type directed manner. The intermediate representation (IR) obtained from erasure carries metadata that mark the linearity of certain critical expressions. For example, metadata is attached to **match**-expressions to indicate whether the scrutinee is linear or unrestricted. This information is used to guide further optimizations for improving runtime efficiency.

One of the optimizations performed is constructor unboxing. The layouts of inductive type constructors are analyzed to determine if the inductive type is suitable for unboxing. For example, consider the singleton type defined as follows:

inductive sing
$$\langle s \rangle \% \{A : \mathsf{Type} \langle s \rangle \} (x : A) := \mathsf{Just} : \forall (x : A) \to \mathsf{sing} \ x$$

Here, Just is the only constructor of type sing. This means that pattern matching on a value of type sing is redundant as there is only one possible case. Expressions of the form Just m are unboxed to m to reduce the number of indirections at runtime. In general, an inductive type can be unboxed if it has a single constructor and the constructor has a single non-ghost field.

To reduce the time spent on allocating and deallocating heap objects, we utilize in-place updates for linear values. This optimization is similar to recent works on function in-place programming [27, 34] where allocated heap memory is reused instead of being garbage collected. Unlike these works which utilize reference counting to dynamically check the viability of an in-place update, the metadata in our IR is sufficient to statically determine if an in-place optimization is safe.

7 Related Work

Session types are a class of type systems pioneered by Honda [23] for structuring dyadic communication in the π -calculus. Abramsky notices deep connections between the Linear Logic [21] of Girard and concurrency, predicting that Linear Logic will play a foundation role in future theories of concurrent computation [2, 3]. Caires and Pfenning show an elegant correspondence between session types and Linear Logic [12]. Gay and Vasconcelos integrate session types with λ -calculus [19] which allows one to express concurrent processes using standard functional programming. Wadler further refines the calculus of Gay and Vasconcelos to be deadlock free by construction [45].

Toninho together with Caires and Pfenning develops the first dependent session type systems [33, 39]. These works extend the existing logic of Caires and Pfenning [12] with universal and existential quantifiers to precisely specify properties of communicated messages.

Toninho and Yoshida present an interesting language [40] that integrates both π -calculus style processes and λ -calculus style terms using a contextual monad. Additionally, full λ -calculi are embedded in both functional types and session types to enable large elimination.

Wu and Xi [46] implement session types in the ATS programming language [48] which supports DML style dependent types [47]. This allows them to specify the properties of concurrent programs and verify them using proof automation. While DML style dependency is well suited for automatic reasoning, certain properties can be difficult to encode due to restrictions on the type level language.

Thiemann and Vasconcelos [38] introduce the LDST calculus which utilizes label dependent session types to elegantly describe communication patterns. Communication protocols written in non-dependent session type systems can essentially be simulated through label dependency. On

the other hand, LDST's minimalist design limits its capabilities for general verification as label dependency by itself is too weak to express many interesting program properties.

Das and Pfenning develop a refinement session type system [16] where the types of concurrent programs can be refined with logical predicates. Similarly to DML style dependent types, the expressiveness of refinement session types is intentionally limited to facilitate proof automation. The Martin-Löf style dependent session types of TLL_C allow users to express and verify more complex program properties at the cost of decidable proof automation.

Atkey proposes QTT [5] based on initial ideas of McBride [30]. QTT is a dependent type theory which tracks resource usage through semi-ring annotations on binders. By instantiating the semi-ring and its ordering relation correctly, QTT can simulate linear types. The Idris 2 programming language [11] (based on QTT) implements a session typed DSL [10] around its raw communication primitives. The authors do not formalize these session types or study its meta-theory. Unlike TLL_C where a library provider could specify a type (such as channels) as linear and automatically enforce its usage in client code through type checking, the obligation of resource tracking is pushed to the client in QTT where binders must be correctly annotated a priori. User mistakes in the annotations could lead to resources being improperly tracked in a program despite passing type checking.

Hinrichsen et al. develop Actris [22] which extends the Iris [25] separation logic framework with dependent separation protocols. Compared to our work, Actris reasons about concurrent programs at a lower level of abstraction. This gives it greater precision and flexibility when dealing with imperative and unsafe programming features. However, the low level nature of Actris reduces its effectiveness at providing guidance for writing programs. In this regard, the interactivity of type systems is more beneficial to helping users construct correct programs in the first place.

8 Conclusion

 TLL_C is a linear dependently typed programming language which extends the TLL type theory with dependent session types. Through examples, we demonstrate how dependent session types can be effectively applied to verify concurrent programs. The expressive power of Martin-Löf style dependency allows TLL_C session types to capture the expected semantics of concurrent programs. This results in greater verification precision and flexibility when compared to other type systems with more restricted forms of dependency. We study the meta-theory of TLL_C and show that it is sound as both a term calculus and also as a process calculus. A prototype compiler is implemented which compiles TLL_C programs into safe concurrent C code.

A direction of research we intend to explore is the integration of dependency with multi-party session types [24]. Protocols expressed through such a session type system will be able to coordinate interactions between processes from a global viewpoint. We predict dependency will again play a key role in verifying the correctness of multi-party concurrent computation.

References

- [1] Andreas Abel and Brigitte Pientka. 2011. Higher-Order Dynamic Pattern Unification for Dependent Types and Records. In *Typed Lambda Calculi and Applications*, Luke Ong (Ed.). Vol. 6690. Springer Berlin Heidelberg, Berlin, Heidelberg, 10–26. doi:10.1007/978-3-642-21691-6 5 Series Title: Lecture Notes in Computer Science.
- [2] Samson Abramsky. 1993. Computational interpretations of linear logic. Theoretical Computer Science 111, 1 (1993), 3-57. doi:10.1016/0304-3975(93)90181-R
- [3] Samson Abramsky. 1994. Proofs as processes. Theoretical Computer Science 135, 1 (1994), 5-9. doi:10.1016/0304-3975(94)00103-0
- [4] Agda development team. 2023. Agda 2.6.4 documentation. https://agda.readthedocs.io/en/v2.6.4
- [5] Robert Atkey. 2018. The Syntax and Semantics of Quantitative Type Theory. In LICS '18: 33rd Annual ACM/IEEE Symposium on Logic in Computer Science, July 9–12, 2018, Oxford, United Kingdom. doi:10.1145/3209108.3209189
- [6] Stephanie Balzer and Frank Pfenning. 2017. Manifest sharing with session types. *Proceedings of the ACM on Programming Languages* 1, ICFP (Aug. 2017), 1–29. doi:10.1145/3110281

[7] Bruno Barras and Bruno Bernardo. 2008. The Implicit Calculus of Constructions as a Programming Language with Dependent Types. In Foundations of Software Science and Computational Structures, Roberto Amadio (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 365–379.

- [8] Giovanni Bernardi and Matthew Hennessy. 2016. Using higher-order contracts to model session types. Logical Methods in Computer Science 12 (06 2016). doi:10.2168/LMCS-12(2:10)2016
- [9] Giovanni Tito Bernardi, Ornela Dardha, Simon J. Gay, and Dimitrios Kouzapas. 2014. On Duality Relations for Session Types. In *TGC*.
- [10] Edwin Brady. 2021. Idris 2: Quantitative Type Theory in Practice. In 35th European Conference on Object-Oriented Programming (ECOOP 2021) (Leibniz International Proceedings in Informatics (LIPIcs), Vol. 194), Anders Møller and Manu Sridharan (Eds.). Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany, 9:1–9:26. doi:10.4230/ LIPIcs.ECOOP.2021.9
- [11] Edwin C. Brady. 2021. Idris 2: Quantitative Type Theory in Practice. CoRR abs/2104.00480 (2021). arXiv:2104.00480 https://arxiv.org/abs/2104.00480
- [12] Luís Caires and Frank Pfenning. 2010. Session Types as Intuitionistic Linear Propositions. 222–236. doi:10.1007/978-3-642-15375-4 16
- [13] Iliano Cervesato and Frank Pfenning. 2002. A Linear Logical Framework. Information and Computation 179, 1 (2002), 19–75. doi:10.1006/inco.2001.2951
- [14] Adam Chlipala. 2013. Certified Programming with Dependent Types: A Pragmatic Introduction to the Coq Proof Assistant. The MIT Press.
- [15] Jesper Cockx and Andreas Abel. 2018. Elaborating Dependent (Co)Pattern Matching. Proc. ACM Program. Lang. 2, ICFP, Article 75 (jul 2018), 30 pages. doi:10.1145/3236770
- [16] Ankush Das and Frank Pfenning. 2020. Verified Linear Session-Typed Concurrent Programming. In Proceedings of the 22nd International Symposium on Principles and Practice of Declarative Programming (Bologna, Italy) (PPDP '20). Association for Computing Machinery, New York, NY, USA, Article 7, 15 pages. doi:10.1145/3414080.3414087
- [17] W. Diffie and M. Hellman. 1976. New directions in cryptography. IEEE Transactions on Information Theory 22, 6 (1976), 644–654. doi:10.1109/TIT.1976.1055638
- [18] Qiancheng Fu and Hongwei Xi. 2023. A Two-Level Linear Dependent Type Theory. arXiv:2309.08673 [cs.PL]
- [19] Simon Gay and Vasco Vasconcelos. 2010. Linear type theory for asynchronous session types. J. Funct. Program. 20 (01 2010), 19–50. doi:10.1017/S0956796809990268
- [20] Simon J. Gay, Peter Thiemann, and Vasco Thudichum Vasconcelos. 2020. Duality of Session Types: The Final Cut. *ArXiv* abs/2004.01322 (2020), 23–33.
- [21] Jean-Yves Girard. 1987. Linear logic. Theoretical Computer Science 50, 1 (1987), 1-101. doi:10.1016/0304-3975(87)90045-4
- [22] Jonas Kastberg Hinrichsen, Jesper Bengtson, and Robbert Krebbers. 2019. Actris: Session-Type Based Reasoning in Separation Logic. *Proc. ACM Program. Lang.* 4, POPL, Article 6 (dec 2019), 30 pages. doi:10.1145/3371074
- [23] Kohei Honda. 1993. Types for Dyadic Interaction. In CONCUR.
- [24] Kohei Honda, Nobuko Yoshida, and Marco Carbone. 2016. Multiparty Asynchronous Session Types. J. ACM 63, 1, Article 9 (mar 2016), 67 pages. doi:10.1145/2827695
- [25] Ralf Jung, David Swasey, Filip Sieczkowski, Kasper Svendsen, Aaron Turon, Lars Birkedal, and Derek Dreyer. 2015. Iris: Monoids and Invariants as an Orthogonal Basis for Concurrent Reasoning. SIGPLAN Not. 50, 1 (Jan. 2015), 637–650. doi:10.1145/2775051.2676980
- [26] Pierre Letouzey. 2003. A New Extraction for Coq. In Types for Proofs and Programs, Gerhard Goos, Juris Hartmanis, Jan Van Leeuwen, Herman Geuvers, and Freek Wiedijk (Eds.). Vol. 2646. Springer Berlin Heidelberg, Berlin, Heidelberg, 200–219. doi:10.1007/3-540-39185-1_12 Series Title: Lecture Notes in Computer Science.
- [27] Anton Lorenzen, Daan Leijen, and Wouter Swierstra. 2023. FP²: Fully in-Place Functional Programming. Proc. ACM Program. Lang. 7, ICFP, Article 198 (aug 2023), 30 pages. doi:10.1145/3607840
- [28] Zhaohui Luo and Y Zhang. 2016. A Linear Dependent Type Theory. 69-70.
- [29] Per Martin-Löf. 1975. An Intuitionistic Theory of Types: Predicative Part. In Logic Colloquium '73, H.E. Rose and J.C. Shepherdson (Eds.). Studies in Logic and the Foundations of Mathematics, Vol. 80. Elsevier, 73–118. doi:10.1016/S0049-237X(08)71945-1
- [30] Conor McBride. 2016. I Got Plenty o' Nuttin'. In A List of Successes That Can Change the World.
- [31] Peter O'Hearn. 2003. On bunched typing. Journal of Functional Programming 13, 4 (July 2003), 747–796. doi:10.1017/ S0956796802004495
- [32] Peter W. O'Hearn and David J. Pym. 1999. The Logic of Bunched Implications. Bulletin of Symbolic Logic 5, 2 (June 1999), 215–244. doi:10.2307/421090
- [33] Frank Pfenning, Luis Caires, and Bernardo Toninho. 2011. Proof-Carrying Code in a Session-Typed Process Calculus. In *Certified Programs and Proofs*, Jean-Pierre Jouannaud and Zhong Shao (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 21–36.

- [34] Alex Reinking*, Ningning Xie*, Leonardo de Moura, and Daan Leijen. 2020. Perceus: Garbage Free Reference Counting with Reuse (Extended version). Technical Report MSR-TR-2020-42. Microsoft. https://www.microsoft.com/en-us/research/publication/perceus-garbage-free-reference-counting-with-reuse/ (*) The first two authors contributed equally to this work. v4, 2021-06-07. Extended version of the PLDI'21 paper..
- [35] Matthieu Sozeau, Simon Boulier, Yannick Forster, Nicolas Tabareau, and Théo Winterhalter. 2020. Coq Coq correct! verification of type checking and erasure for Coq, in Coq. Proceedings of the ACM on Programming Languages 4, POPL (Jan. 2020), 1–28. doi:10.1145/3371076
- [36] The Coq Development Team. 2020. The Coq Proof Assistant, version 8.11.0. doi:10.5281/ZENODO.3744225
- [37] The Rust teams. 2022. Rust Programming Language. http://www.rust-lang.org/
- [38] Peter Thiemann and Vasco T. Vasconcelos. 2019. Label-Dependent Session Types. Proc. ACM Program. Lang. 4, POPL, Article 67 (dec 2019), 29 pages. doi:10.1145/3371135
- [39] Bernardo Toninho, Luís Caires, and Frank Pfenning. 2011. Dependent Session Types via Intuitionistic Linear Type Theory. In Proceedings of the 13th International ACM SIGPLAN Symposium on Principles and Practices of Declarative Programming (Odense, Denmark) (PPDP '11). Association for Computing Machinery, New York, NY, USA, 161–172. doi:10.1145/2003476.2003499
- [40] Bernardo Toninho and Nobuko Yoshida. 2018. Depending on Session-Typed Processes. 128–145. doi:10.1007/978-3-319-89366-2_7
- [41] David N. Turner and Philip Wadler. 1999. Operational interpretations of linear logic. *Theoretical Computer Science* 227, 1 (1999), 231–248. doi:10.1016/S0304-3975(99)00054-7
- [42] Matthijs Vákár. 2014. Syntax and Semantics of Linear Dependent Types. CoRR abs/1405.0033 (2014). arXiv:1405.0033 http://arxiv.org/abs/1405.0033
- [43] P. Wadler. 1990. Linear Types can Change the World!. In Programming Concepts and Methods.
- [44] Philip Wadler. 1991. Is There a Use for Linear Logic? SIGPLAN Not. 26, 9 (May 1991), 255-273. doi:10.1145/115866.115894
- [45] Philip Wadler. 2012. Propositions as Sessions. In Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming (Copenhagen, Denmark) (ICFP '12). Association for Computing Machinery, New York, NY, USA, 273–286. doi:10.1145/2364527.2364568
- [46] Hanwen Wu and Hongwei Xi. 2017. Dependent Session Types. CoRR abs/1704.07004 (2017). arXiv:1704.07004 http://arxiv.org/abs/1704.07004
- [47] Hongwei Xi. 2007. Dependent ML An approach to practical programming with dependent types. Journal of Functional Programming 17, 2 (2007), 215–286. doi:10.1017/S0956796806006216
- [48] Hongwei Xi. 2010. The ATS Programming Language. http://www.ats-lang.org/