

# CS442

## Module 9: Concurrent Programming

University of Waterloo

Winter 2023

“Concurrency is the most extreme form of programming. It’s like white-water rafting without the raft or sky diving without the parachute.”

— Peter Buhr

### 1 Concurrent Programming

We’ve developed two core calculi and looked at four exemplar programming languages, but we’re as yet missing one of the most important characteristics of modern computing: *concurrency*. First, let’s clarify some terms, because the concurrent programming community are quite picky about these terms:

- *Parallelism* is the *phenomenon* of two or more things—presumably, two computations—actually happening at the same time.
- *Concurrency* is the *experience* of two or more things—such as applications, tasks, or threads—appearing to happen at the same time, whether or not they actually are.

For instance, an average computer in the 1990s had only one CPU and one core, and so had no parallelism, but could still run multiple applications “at the same time”, and so still had concurrency. In this case, concurrency is achieved by quickly switching which task the CPU is concerned with between the various running programs. Most home computers of the 1980s had neither: we could only run one program at a time, and we liked it; it builds character! It is also possible for a system to be parallel without being concurrent: for instance, a compiler may optimize an imperative loop into a specialized parallel operation on a particular CPU, but this is only visible to the programmer as a speed boost, so the programmer’s experience still lacks concurrency. And, of course, a system can have both: on a modern, multi-core CPU, one uses concurrency to take advantage of the parallelism, by running multiple programs or threads on their multiple cores. The concurrent tasks appear to happen at the same time because they *do* happen at the same time.

The domain of concurrent programming has changed dramatically over time, in response to growing parallelism. Early exploration into concurrency was theoretical, and then became practical with networks. When multi-core CPUs started becoming the norm, concurrency rose in importance again, because it is often impossible to take advantage of real parallelism without concurrency.

The core idea of a concurrent system is that there are multiple tasks which can communicate with each other, and computation can proceed in any of them at any rate. This is contrary to the behavior of the  $\lambda$ -calculus or the Simple Imperative Language, and contrary to all of our exemplar languages (at least, as far as we’ve investigated them), so we will need both a new formal language and a new exemplar language for concurrency. In practice, of course, in the same way that many languages which are not fundamentally object oriented have picked up object-oriented features, most languages which are not fundamentally concurrent have picked up at least some kind of concurrent features.

In a formal model of concurrency, we need a way of expressing multiple simultaneous tasks, and, as usual, a way of taking a step. Unlike our previous formal semantics, we *want* this semantics to be non-deterministic: the fact

that any of multiple tasks may proceed is the essence of concurrency. However, we will not define our semantics as truly parallel: a step will involve only one task taking a step, with the non-determinism only affecting *which* task that is.

Another issue to be addressed is *models* of concurrency, i.e., how a concurrent language presents multiple tasks. Models of concurrency are defined across two dimensions: how one *specifies* multiple tasks, and how those multiple tasks *communicate*. Specification comes down to what structures a language provides for a program to create tasks, and often, to specify real parallelism as well. Options include threads of various sorts, actors, processes, and many others, but all of these terms are imprecise and ambiguous. We will put little focus on the specification of concurrency beyond how it's done in our formal model and exemplar language, and instead will focus, as usual, on semantics.

The other dimension is communication, which falls largely into two forms: *shared-memory* concurrency and *message-passing* concurrency. In shared-memory concurrency, either the entire heap or some portion of the heap is shared between multiple tasks. In our terms, all tasks have the same  $\Sigma$ , and if they share any labels, they can see changes made by other tasks. However, because multiple tasks can compute at the same time, there may not be any guarantee that one task's write to a location in  $\Sigma$  occurs before another task's read. Other mechanisms, such as locks, are required to guarantee ordering.

In message-passing concurrency, two new primitive operations are introduced: *sending* a message and *waiting* for a message. A task which is waiting for a message cannot proceed until another task sends it a message. A task may send a message at any time, but to do so, must have a way of communicating with the target task. Thus, message-passing concurrency requires some form of *message channels*, by which two tasks can arrange to exchange messages. Ordering is guaranteed by the directionality of message passing; a waiting task will not proceed until a sending task sends a message.

**Aside:** Shared-memory and message-passing concurrency are equally powerful, and in fact, either can be rewritten in terms of the other. However, this rewriting is fairly unsatisfying: shared memory can be rewritten in terms of message passing by imagining  $\Sigma$  itself as a “task” that expects messages instructing it to read and write certain memory locations. Equivalently, we can consider a modern CPU as sending messages to the memory bus, rather than simply reading and writing memory. As a practical matter, message-passing implementations tend to be slower than shared-memory implementations of the same algorithm.

This course is not intended to compete with CS343 (Concurrent and Parallel Programming), so will take a very language-focused view of concurrency. That course uses concurrency to solve problems; in this course, concurrency is only a cause of problems.

## 2 $\pi$ -Calculus

In the  $\lambda$ -calculus, we built a surprising amount of functionality around abstractions: with only three constructs in the language (abstractions, applications, and variables), we could represent numbers, conditionals, and ultimately, anything computable.  $\pi$ -calculus (The Pi Calculus) takes a similar approach to message-passing concurrency. The only structures are concurrency, communication, replication, and names, but these will be sufficient to build any computation. Notably lacking is functions (or abstractions).  $\pi$ -calculus itself was developed by Robin Milner<sup>1</sup>, Joachim Parrow, and David Walker in *A Calculus of Mobile Processes*[1], but it was the culmination of a long line of development of calculi of communicating processes, to which Uffe Engberg and Mogens Nielsen also made significant contributions.

We will look at the behavior of  $\pi$ -calculus, but will not discuss how to encode complex computation into  $\pi$ -calculus concurrency.  $\pi$ -calculus is Turing-complete, but like the  $\lambda$ -calculus, it's more common to layer other semantics on top of it than to take advantage of its own computational power.

A  $\pi$ -calculus program consists of any number of concurrent tasks, called *processes*, separated by pipes (`|`). Those tasks can be grouped by parentheses. Each process can *create* a channel, *receive* a message on a channel, *send* a message on a channel, *replicate* processes, or *terminate*. The only values in  $\pi$ -calculus are channels, so any encoding

<sup>1</sup>The same Milner of Hindley-Milner type inference.

of useful information also needs to be done with channels, and the only value you can send over a channel is another channel. Channels are also called “names”, because they are simply named by variables; two processes must agree on a name in order to communicate (as well as another restriction which we will discuss soon).

The syntax of  $\pi$ -calculus is as follows, presented in BNF, with  $\langle program \rangle$  as the starting non-terminal:

$$\begin{aligned}
 \langle program \rangle &::= \langle program \rangle \text{ “|” } \langle program \rangle \\
 &| \langle receive \rangle \\
 &| \langle send \rangle \\
 &| \langle restrict \rangle \\
 &| \langle replicate \rangle \\
 &| \langle terminate \rangle \\
 \langle receive \rangle &::= \langle var \rangle (\langle var \rangle). \langle program \rangle \\
 \langle send \rangle &::= \overline{\langle var \rangle} \text{ “(” } \langle var \rangle \text{ “)” } . \langle program \rangle \\
 \langle restrict \rangle &::= (\nu \langle var \rangle) \langle program \rangle \\
 \langle replicate \rangle &::= ! \langle program \rangle \\
 \langle terminate \rangle &::= 0 \\
 \langle var \rangle &::= a|b|c|\dots
 \end{aligned}$$

Note that  $\nu$  is the Greek letter nu, not the Latin/English letter ‘v’, because somebody decided that using confusing, ambiguous Greek letters was acceptable; we will avoid using  $v$  as a variable for this reason. Like in the  $\lambda$ -calculus, we will actually be more lax in our use of variable names than this BNF suggests, for clarity. Like in the Simple Imperative Language, this is assumed to be an abstract syntax, and we will add parentheses to disambiguate as necessary. The pipe (|) in  $\langle program \rangle$  has the lowest precedence, so, for instance,  $x(y).0|z(a).0$  is read as  $(x(y).0)|(z(a).0)$ , not  $x(y).(0|z(a).0)$ .

Unfortunately,  $\pi$ -calculus uses several of the symbols which we use in BNF as well, which we’ve surrounded in quotes to separate them from BNF.  $\pi$ -calculus also uses an overline, which is shown in the BNF above, but is at the very least unusual. Here are two small examples to clarify the syntax. The following snippet receives a message on the channel  $x$ , into the variable  $y$ , before proceeding with the process  $P$ :

$$x(y).P$$

The following snippet sends  $y$  on the channel  $x$ , before proceeding with the process  $P$ :

$$\overline{x} \langle y \rangle .P$$

As discussed, a program consists of a number of processes, separated by pipes. Each process is itself a program, so the distinction is just usage. We will use the term “process” to refer to any construction *other* than the composition of multiple programs with a pipe, so that a program can be read as a list of processes.

A program proceeds through its processes sending and receiving messages on channels until they terminate. For instance, this program consists of two processes, of which the first sends the message  $h$  to the second, and the second then attempts to pass that  $h$  along on another channel:

$$\overline{x} \langle h \rangle .0|x(y).\overline{z} \langle y \rangle .0$$

The first process is  $\overline{x} \langle h \rangle .0$ , which consists of a send of  $h$  over the channel  $x$ , and then termination of the process (0). The second process is  $x(y).\overline{z} \langle y \rangle .0$ , which consists of a receive of  $y$  from the channel  $x$ , then a send of  $y$  over the channel  $z$ , then termination. Programs in  $\pi$ -calculus proceed by sending and receiving messages; in this case, the program can proceed, because a process is trying to send on  $x$ , and another process is trying to receive on  $x$ . After sending that message, the program looks like this:

$$0|\overline{z} \langle h \rangle .0$$

Like in  $\lambda$ -calculus applications, message receipt works by substitution, so the  $y$  was substituted for the received message,  $h$ . We will not formally define substitution for  $\pi$ -calculus; it is fairly straightforward. A terminating process can simply be removed, so the next step is as follows:

$$\bar{z} \langle h \rangle . 0$$

This program cannot proceed, because no process is prepared to receive a message on the channel  $z$ .

Consider this similar program:

$$\bar{x} \langle h \rangle . z(a). 0 | x(y). \bar{z} \langle y \rangle . 0 | x(y). 0$$

This time, we have three processes. The first sends the message  $h$  over the channel  $x$ , then receives a message on the channel  $z$ , then terminates. The second is identical to our original second process: it receives a message on  $x$ , then sends it back on  $z$ . The third message receives a message on  $x$ , then immediately terminates. This program can proceed, because a process is sending on  $x$  and a process is prepared to receive on  $x$ . But, which process receives the message?  $\pi$ -calculus is non-deterministic, so the answer is that *either* process may receive the message. Both ways for the program to proceed are valid. In this case, these two reductions are both valid:

$$\begin{aligned} \Rightarrow z(a). 0 | \bar{z} \langle h \rangle . 0 | x(y). 0 & \Rightarrow z(a). 0 | x(y). \bar{z} \langle y \rangle . 0 | 0 \\ \Rightarrow 0 | 0 | x(y). 0 & \Rightarrow z(a). 0 | x(y). \bar{z} \langle y \rangle . 0 \\ \Rightarrow 0 | x(y). 0 & \\ \Rightarrow x(y). 0 & \end{aligned}$$

The first sequence of reductions occurs if the message is received by the second process, and the third occurs if the message is received by the third process. The first sequence may seem more complete, since two of the three processes terminated, but both are valid. We are using  $\Rightarrow$  informally for “takes a step”, because we haven’t yet formally defined  $\rightarrow$ , and we will discover that there’s an extra complication to the definition of steps in  $\pi$ -calculus in Section 2.2.

Because name uniqueness is so important to communication,  $\pi$ -calculus also has a mechanism to *restrict* names. For instance, consider this rewrite of the above program:

$$(\nu x)(\bar{x} \langle h \rangle . z(a). 0 | x(y). \bar{z} \langle y \rangle . 0) | x(y). 0$$

It is identical to the previous program, except that the first two processes are nested inside of a *restriction*:  $(\nu x)$ . A restriction is like a variable binding, in that it scopes the variable to its context: the  $x$  inside of the restriction is not the same as the  $x$  outside of the restriction. Unlike variable binding, however, it doesn’t bind it to any particular value—remember, names are all there are in  $\pi$ -calculus, so there’s nothing else it could be bound *to*—it merely makes for two distinct interpretations of the name. That’s why it’s called a restriction; it restricts the meaning of, in this case,  $x$ , within the restriction expression. Now, this program can only proceed like so:

$$\begin{aligned} \Rightarrow (\nu x)(z(a). 0 | \bar{z} \langle h \rangle . 0) | x(y). 0 \\ \Rightarrow (\nu x)(0 | 0) | x(y). 0 \\ \Rightarrow (\nu x)(0) | x(y). 0 \\ \Rightarrow x(y). 0 \end{aligned}$$

In the last step, we can remove a restriction when it is no longer restricting anything (i.e., when its contained process terminates). When a restriction is at the top level like this, it can always be rewritten by renaming variables to new, fresh names, so the above program is equivalent to the following program:

$$\bar{x'} \langle h \rangle . z(a). 0 | x'(y). \bar{z} \langle y \rangle . 0 | x(y). 0$$

**Aside:** If your eyes are glazing over from  $\pi$ -calculus syntax, don’t worry, you’re not the only one. Something about  $\pi$ -calculus’s use of overlines and  $\nu$  and pipes makes it semantically dense and difficult to read. I have no advice to alleviate this; just be careful of the pipes and parentheses.

A restriction only applies to the variable it names, so processes within restrictions are still allowed to communicate with processes outside of restrictions:

$$\begin{aligned} & (\nu x)\bar{z}\langle h \rangle .0|z(y).\bar{x}\langle y \rangle .0 \\ \Rightarrow & (\nu x)0|\bar{x}\langle h \rangle .0 \\ \Rightarrow & \bar{x}\langle h \rangle .0 \end{aligned}$$

Substitution must be aware of restriction, because a restricted variable is distinct from the same name in the surrounding code. For instance:

$$(x(y).(\nu x).x(y).0)[z/x] = (z(y).(\nu x).x(y).0)$$

This exception is the same as is introduced by  $\lambda$ -abstractions in substitution for the  $\lambda$ -calculus.

The only messages that processes can send are names. Names are also the channels by which processes send messages. As a consequence, processes can send channels over channels. For instance, consider this program:

$$\bar{x}\langle z \rangle .0|x(y).\bar{y}\langle h \rangle .0|z(a).0$$

The first process will send the name  $z$  over the channel  $x$ . The second process is waiting to receive a message on the channel  $x$ . The third process is waiting to receive a message on the channel  $z$ , but there is no send on the channel  $z$  in the entire program. The third process cannot possibly proceed, but the first and second can, like so:

$$\Rightarrow 0|\bar{z}\langle h \rangle .0|z(a).0 \Rightarrow \bar{z}\langle h \rangle .0|z(a).0$$

The second process received a  $z$  on the channel  $x$ , as the variable  $y$ . But, it then proceeds to *send* on the channel  $y$ . Because message receipt works by substitution, that  $y$  has been substituted for  $z$ . This program can now proceed, by sending  $h$  to the third process (at which point both the second and third process terminate).

Restriction has an unusual interaction with sending channels. For instance, consider this program:

$$(\nu x)\bar{z}\langle x \rangle .x(y).0|z(a).\bar{a}\langle x \rangle .0$$

The first process is under a restriction for  $x$ , and the second process is not. But, the behavior of the first process is to send  $x$  over the channel  $z$ , and the second process is waiting to receive a message on the channel  $z$ . It's then going to send  $x$  back, but that  $x$  isn't the same  $x$  as the first process's  $x$ , because of the restriction. So, what happens if we send a restricted channel outside of its own restriction? How do we deal with these two conflicting  $x$ 's? The answer is made clear by our statement that a restriction can always be rewritten by simply using a new name. In this case, this program can be rewritten like so:

$$\bar{z}\langle x' \rangle .x'(y).0|z(a).\bar{a}\langle x \rangle .0$$

From this state, the steps are clear:

$$\begin{aligned} & \Rightarrow x'(y)0|\bar{x'}\langle x \rangle .0 \\ & \Rightarrow 0|0 \\ & \Rightarrow 0 \end{aligned}$$

Finally,  $\pi$ -calculus supports process creation: a process may create more processes. There are actually two mechanisms of process creation. First, processes may simply be nested. For instance, consider this program:

$$x(y).(\bar{y}\langle h \rangle .0|\bar{y}\langle m \rangle .0)|f(a).\bar{z}\langle a \rangle .0|f(b).\bar{z}\langle b \rangle .0|\bar{x}\langle f \rangle .0$$

Note in particular the position of the parentheses: this program has *four* processes, not five! The first process is  $x(y).(\bar{y}\langle h \rangle .0|\bar{y}\langle k \rangle .0)$ . Although this process has the pipe which separates multiple processes within it, those are not two independent processes until this process has received a message on the channel  $x$ . Essentially, as soon as this process receives a message, it will split into two processes. This program can proceed as follows (this is not the only possible sequence):

$$\begin{aligned}
&\Rightarrow (\bar{f} \langle h \rangle . 0 | \bar{f} \langle m \rangle . 0 | f(a) . \bar{z} \langle a \rangle . 0 | f(b) . \bar{z} \langle b \rangle . 0 \\
&\Rightarrow \bar{f} \langle h \rangle . 0 | \bar{f} \langle m \rangle . 0 | f(a) . \bar{z} \langle a \rangle . 0 | f(b) . \bar{z} \langle b \rangle . 0 \\
&\Rightarrow 0 | \bar{f} \langle m \rangle . 0 | \bar{z} \langle h \rangle . 0 | f(b) . \bar{z} \langle b \rangle . 0 \\
&\Rightarrow \bar{f} \langle m \rangle . 0 | \bar{z} \langle h \rangle . 0 | f(b) . \bar{z} \langle b \rangle . 0 \\
&\Rightarrow 0 | \bar{z} \langle h \rangle . 0 | \bar{z} \langle m \rangle . 0 \\
&\Rightarrow \bar{z} \langle h \rangle . 0 | \bar{z} \langle m \rangle . 0
\end{aligned}$$

**Exercise 1.** Give another possible sequence for this example.

The other mechanism of process creation is *replication*. The ! operator creates an endless sequence of identical processes. For instance, consider the following program:

$$!x(a).0|\bar{x}\langle b\rangle.0|\bar{x}\langle c\rangle.0|\bar{x}\langle d\rangle.0$$

There are three processes trying to send on the channel  $x$ , but only one process with a receive on the channel  $x$ . However, the ! creates any number of the same process, so all three of the sending processes can proceed, in any order. This is one possible sequence:

$$\begin{aligned}
&\Rightarrow x(a).0|!x(a).0|\bar{x}\langle b\rangle.0|\bar{x}\langle c\rangle.0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow 0|!x(a).0|\bar{x}\langle b\rangle.0|0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow !x(a).0|\bar{x}\langle b\rangle.0|0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow !x(a).0|\bar{x}\langle b\rangle.0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow x(a).0|!x(a).0|\bar{x}\langle b\rangle.0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow 0|!x(a).0|0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow !x(a).0|0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow !x(a).0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow x(a).0|!x(a).0|\bar{x}\langle d\rangle.0 \\
&\Rightarrow 0|!x(a).0|0 \\
&\Rightarrow !x(a).0|0 \\
&\Rightarrow !x(a).0
\end{aligned}$$

## 2.1 Concurrency vs. Parallelism

The astute reader may have noticed that we have described sequences of steps, with no true parallelism. For instance, consider the following program:

$$\bar{x}\langle a \rangle . 0 | \bar{y}\langle b \rangle . 0 | x(z) . 0 | y(z) . 0$$

There are two possible steps this program can take—it can send a message on  $x$  or  $y$ —but we describe it as taking one or the other, not both at the same time. The concurrency comes from the lack of prioritization, and non-determinism: each of these two options is equally valid, and to consider how this program proceeds, we need to consider both possibilities. But, the concurrency is restricted by the nature of messages: only pairs of matching sends and receives can actually proceed.

Because concurrency models the *appearance* of multiple tasks happening simultaneously, in most cases, it is not necessary to model true parallelism. The most complex<sup>2</sup> formal models in the domain of concurrency and parallelism are models of parallel shared-memory architectures, and even they are formally descriptions of concurrency rather than parallelism, in that they model parallel action as a non-deterministic ordering.

---

<sup>2</sup>In this author's opinion.

## 2.2 Structural Congruence

Because processes may proceed in any order in  $\pi$ -calculus,  $P|Q$  is not meaningfully distinct from  $Q|P$ . Similarly,  $(\nu x)P|Q$  is not meaningfully distinct from  $P[y/x]|Q$ , where  $y$  is a new name, and  $\alpha$ -equivalent programs are also indistinct.

In the  $\lambda$ -calculus, these equivalences mostly gave us a baseline for comparing things. In  $\pi$ -calculus, it would be difficult or impossible to define reduction without this equivalence, because of the non-deterministic ordering of steps.

This equivalence is defined formally as *structural congruence*, written as  $\equiv$ . That is,  $P \equiv Q$  means that  $P$  is structurally congruent to  $Q$ . Structural congruence is reflexive, symmetric, and transitive.

The formal rules for structural congruence follow. Note that different presentations of  $\pi$ -calculus present slightly different but equivalent rules of structural congruence, so this may not exactly match other materials on the same topic.

### Definition 1. (Structural congruence)

Let the metavariables  $P$ ,  $Q$ , and  $R$  range over programs, and  $x$  and  $y$  range over names. Then the following rules describe structural congruence of  $\pi$ -calculus programs:

$$\begin{array}{lll}
 \text{C\_ALPHA} \quad \frac{P =_{\alpha} Q}{P \equiv Q} & \text{C\_ORDER} \quad \frac{}{P|Q \equiv Q|P} & \text{C\_NEST} \quad \frac{P \equiv P'}{P|Q \equiv P'|Q} \\
 \\
 \text{C\_PAREN} \quad \frac{}{(P|Q)|R \equiv P|(Q|R) \equiv P|Q|R} & \text{C\_TERMINATION} \quad \frac{}{0|P \equiv P} & \\
 \\
 \text{C\_RESTRICTION} \quad \frac{y \text{ is a fresh variable}}{(\nu x)P \equiv P[y/x]} & \text{C\_REPLICATION} \quad \frac{}{!P \equiv P|!P} & 
 \end{array}$$

The C\_ALPHA rule specifies that  $\alpha$ -equivalence implies structural congruence, i.e., two  $\alpha$ -equivalent programs are also structurally congruent. The C\_ORDER and C\_NEST rules allow us to reorder programs and apply structural congruence to subprograms. The C\_PAREN rule specifies that all concurrent processes are equivalent, and different placements of parentheses do not affect their composition, so we can remove parentheses at the top level of a program. The C\_TERMINATION rule describes termination in terms of equivalence: rather than termination being a step, we can describe a program with a terminating process as equivalent to a program without that process. The C\_RESTRICTION rule makes explicit our description of restriction as creating a fresh variable. Finally, the C\_REPLICATION rule makes replication a property of structural congruence, rather than a step: a replicating process is simply equivalent to a version with a replica, and thus, by the transitive property, equivalent to a version with any number of replicas.

Because of C\_TERMINATION, C\_RESTRICTION, and C\_REPLICATION, only sending and receiving messages is described as an actual step of computation. Everything else is structural congruence.

Note that C\_RESTRICTION does *not* allow us to remove all restrictions from a program, because restrictions may be nested inside of other constructs, and structural equivalence does not allow us to enter any other constructs. For instance, the program  $x(y).(\nu z).0$  has no structural equivalent (except for  $\alpha$ -renaming), because the restriction of  $z$  is nested inside of a receipt on the channel  $y$ .

## 2.3 Formal Semantics

With structural congruence, we may now describe the formal semantics of  $\pi$ -calculus.

### Definition 2. (Formal semantics of $\pi$ -calculus)

Let the metavariables  $P$ ,  $Q$ , and  $R$  range over programs, and  $x$ ,  $y$ , and  $z$  range over names. Then the following rules describe the formal semantics of  $\pi$ -calculus:

$$\text{CONGRUENCE } \frac{P \equiv Q \quad Q \rightarrow Q' \quad Q' \equiv R}{P \rightarrow R} \quad \text{MESSAGE } \frac{}{\bar{x}(y).P|x(z).Q|R \rightarrow P|Q[y/z]|R}$$

Because of structural congruence, these two rules are all that is needed to define the semantics of  $\pi$ -calculus. By CONGRUENCE,  $P$  can reduce to  $R$  if  $P$  is equivalent to some  $Q$  which can reduce to some  $Q'$ , and  $Q'$  is equivalent to  $R$ . That is, reduction is ambivalent to structural congruence in either its “from” or “to” state. MESSAGE describes the only actual reduction step in our semantics: if there is a process to send a message, and a process to receive a message on the same channel, then we may take a step in both, by removing the send from the sending process, removing the receipt from the receiving process, and substituting the variable in the receiving process with the value sent by the sending process.

MESSAGE itself is deterministic. The non-determinism in  $\pi$ -calculus is introduced by CONGRUENCE. Every program  $P$  has infinitely many structurally congruent equivalents. Some number of those structurally congruent programs are able to take steps with MESSAGE. Each of those is equivalently correct, and none has priority; all are valid reduction steps.

Consider our previous example:

$$!x(a).0|\bar{x}(b).0|\bar{x}(c).0|\bar{x}(d).0$$

We may now define one possible sequence formally, with structural congruence and reduction:

$$\begin{aligned} & !x(a).0|\bar{x}(b).0|\bar{x}(c).0|\bar{x}(d).0 \\ \equiv & x(a).0|!x(a).0|\bar{x}(b).0|\bar{x}(c).0|\bar{x}(d).0 && (\text{C\_REPLICATION}) \\ \equiv & \bar{x}(c).0|x(a).0|!x(a).0|\bar{x}(b).0|\bar{x}(d).0 && (\text{C\_ORDER}) \\ \rightarrow & 0|0|!x(a).0|\bar{x}(b).0|\bar{x}(d).0 && (\text{MESSAGE}) \\ \equiv & !x(a).0|\bar{x}(b).0|\bar{x}(d).0 && (\text{C\_TERMINATION}) \\ \equiv & x(a).0|!x(a).0|\bar{x}(b).0|\bar{x}(d).0 && (\text{C\_REPLICATION}) \\ \equiv & \bar{x}(b).0|x(a).0|!x(a).0|\bar{x}(d).0 && (\text{C\_ORDER}) \\ \rightarrow & 0|0|!x(a).0|\bar{x}(d).0 && (\text{MESSAGE}) \\ \equiv & !x(a).0|\bar{x}(d).0 && (\text{C\_TERMINATION}) \\ \equiv & x(a).0|!x(a).0|\bar{x}(d).0 && (\text{C\_REPLICATION}) \\ \equiv & \bar{x}(d).0|x(a).0|!x(a).0 && (\text{C\_ORDER}) \\ \rightarrow & 0|0|!x(a).0 && (\text{MESSAGE}) \\ \equiv & !x(a).0 && (\text{C\_TERMINATION}) \end{aligned}$$

## 2.4 The Use of $\pi$ -Calculus

In concurrent programming, most problems—at least, most that are unique to concurrent programming, and aren’t just logic bugs—can be simplified to *happens-before* relationships. That is, with multiple processes able to perform tasks concurrently, you want to guarantee that some task happens before some other task. Concurrent systems are modeled in terms of  $\pi$ -calculus to prove these kinds of happens-before relationships.



For instance, let's say we want to verify that a given program always sends a message on channel  $x$  before sending a message on channel  $y$ . Here is a program that fails to guarantee such a relationship:

$$\bar{a} \langle x \rangle . \bar{b} \langle y \rangle . 0 | a(m) . \bar{m} \langle h \rangle . 0 | b(n) . \bar{n} \langle h \rangle . 0 | x(q) . 0 | y(q) . 0$$

We can demonstrate this by showing a reduction that sends on  $y$  before sending on  $x$ :

$$\begin{aligned} &\rightarrow \bar{b} \langle y \rangle . 0 | \bar{x} \langle h \rangle . 0 | b(n) . \bar{n} \langle h \rangle . 0 | x(q) . 0 | y(q) . 0 \\ &\equiv \bar{b} \langle y \rangle . 0 | b(n) . \bar{n} \langle h \rangle . 0 | \bar{x} \langle h \rangle . 0 | x(q) . 0 | y(q) . 0 \\ &\rightarrow 0 | \bar{y} \langle h \rangle . 0 | \bar{x} \langle h \rangle . 0 | x(q) . 0 | y(q) . 0 \\ &\equiv \bar{y} \langle h \rangle . 0 | y(q) . 0 | \bar{x} \langle h \rangle . 0 | x(q) . 0 \\ &\rightarrow 0 | 0 | \bar{x} \langle h \rangle . 0 | x(q) . 0 \\ &\quad \text{(Premise violated)} \end{aligned}$$

Proving that happens-before relationships hold is, of course, far more complicated, since it is impossible to enumerate the infinitely many possible structurally congruent programs. Luckily, C\_REPLICATION is the only case that can introduce infinite reducible programs, and the difference between them is uninteresting (only how many times the replicated subprogram has been expanded). So, in many cases, it is possible to enumerate all *interesting* reductions. If the program can reduce forever, then it is instead necessary to use inductive proofs for most interesting properties.

Generally,  $\pi$ -calculus is extended with other features to represent the actual computation that each process performs, rather than performing computation through message passing. For instance, the  $\lambda$ -calculus and  $\pi$ -calculus can be overlain directly by allowing processes which contain  $\lambda$ -applications to proceed as in the  $\lambda$ -calculus, while process pairs containing a matching send and receive can proceed as in  $\pi$ -calculus. Such combinations are often used for proving type soundness of concurrent languages.

### 3 Exemplar: Erlang

In  $\pi$ -calculus, we have found a formal semantics for message-passing concurrency. Although there are many programming languages with support for concurrency, and even many programming languages with support for message-passing concurrency, we're looking for an *exemplar* of message-passing concurrency, not just a language with it as a feature. There is a stand-out example which is to message-passing concurrency as Smalltalk is to object orientation: Erlang.

Erlang<sup>3</sup> is a language built on the principle that “everything is a process”. It was created in the late 1980s at Ericsson by Joe Armstrong, Robert Virding, and Mike Williams, to manage telecom systems. There were three primary goals in that context:

- that the system scale from single systems (where many processes would run on one computer) to distributed systems (where processes could be distributed across many computers) with little or no rewriting,
- that processes would be sufficiently isolated that faults in one process would not (necessarily) affect the rest of the system, and
- that individual processes could be replaced live in a running system, allowing for smooth upgrades without any downtime.

These goals led Erlang to a quite extreme design, whereby Erlang programs use processes in the same way as Smalltalk programs use objects. Nearly all compound data is bound in processes, and one interacts with processes by sending and receiving messages. Just like in  $\pi$ -calculus, one can create processes and send channels in messages, allowing sophisticated interactions.

In fact, unlike Smalltalk's objects, Erlang does support some primitive data types which are not processes. Integers, floating point numbers, tuples, lists, key-value maps, and Prolog-like atoms are all supported, and ports—Erlang's name for one end of a communication channel—are not themselves processes (how would one ever send

<sup>3</sup>Pronounced roughly like “air-lang” by people who know how to pronounce words, or “ur-lang” by the kind of troglodytes who pronounce “wiki” as “wick-y”.

a message if message channels were themselves processes which were controlled by messages?). So, it's not quite true that everything is a process, but everything\* is a process. In fact, it's perfectly possible to treat Erlang as a mostly-pure functional language and write totally non-concurrent code. However, we won't focus on that aspect of Erlang. Instead we'll look only at its concurrency features.

Although Erlang has its own unique syntax, by this point, you should be able to guess how most of it works. Like Prolog, variables in Erlang are named with capital letters, and atoms and functions are named with lower-case letters, but its behavior is otherwise more similar to a functional language than a logic language.

### 3.1 Modules

Erlang divides code into modules. We've avoided discussing modules for most of this course, as they're usually uninteresting for language semantics, but creating a process in Erlang requires modules, so we will briefly discuss them. An Erlang file is a module, and must start with a declaration of the module's name. For instance, a module named `sorter` begins as follows:

```
-module(sorter).
```

Most modules define some public functions and some private functions. Any functions which should be usable from other modules must be *exported*. For instance, if we define a function `merge` taking two arguments, we make that function visible like so:

```
-export([merge/2])
```

Note that functions in this context are named with their arity, in this case 2, in the same fashion as Prolog, so `merge/2` is a function named `merge` that takes two arguments.

Functions in the same module can be called with only their name:

```
merge([1, 2, 3], [2, 2, 4])
```

Functions in other modules need to be prefixed with the target module:

```
sorter:merge([1, 2, 3], [2, 2, 4])
```

### 3.2 Processes

A process is created in Erlang with the built-in `spawn` function. `spawn` is called with a module and function name, and the arguments for that function, and the newly created process starts running that function. `spawn` returns a process reference, which can then be used to communicate with the process.

For instance, the following function spawns two processes, passing the process reference of the first to the second. The first process runs the `pong` function in the `pingpong` module with no arguments, and the second runs the `ping` function in the `pingpong` module with the arguments 5 and the reference to the `pong` process:

```
start() ->
  Pong = spawn(pingpong, pong, []),
  spawn(pingpong, ping, [5, Pong]).
```

Generally, functions perform a list of comma-separated actions like this.

Now, let's write the `ping` and `pong` functions. `ping` will send the given number of "ping" messages to the given process, and expect an equal number of "pong" messages in response. `pong` will expect a sequence of "ping" message, and send a "pong" to each. For this to work, we need to know Erlang's syntax for sending and receiving messages.

Messages are sent in Erlang with the `!` operator, as `Target ! Message`. The message can be any Erlang value, but in practice, it is either an atom or a tuple in which the first element is an atom, where the atom specifies the kind of message, and any arguments that the message has fill the rest of the tuple. This convention is simply so that when a message is received, the receiver knows what kind of message it is; the atom is a tag. In our case, the "pong" process does not have a reference to the "ping" process, but the "ping" process does have a reference to the "pong" process, so the "ping" message will need to send a reference along in order for the "pong" process to know how to reply. In  $\pi$ -calculus terms, "ping" must send the channel on which "pong" is to reply.

Messages are received in Erlang with a `receive` expression, which resembles a pattern match, in that it matches the shape of the message received. For a message to be successfully sent, the target process must be running a `receive` with a matching pattern, in the same way that for a  $\pi$ -calculus program to make progress, a sending process must have a matching receiving process. Because `receive` matches particular shapes of messages, a process can receive messages in any order, but process them in the order it chooses, simply by performing `receives` in sequence that match only the kinds of messages it wishes to process.

Knowing this, we will write `ping` first:

```

1 ping(0, _) -> io:format("Ping finished~n", []);
2
3 ping(N, Pong) ->
4     Pong ! {ping, self()},
5     receive
6         pong -> io:format("Pong received~n", [])
7     end,
8     ping(N - 1, Pong).

```

Like in Haskell, functions can be declared in multiple parts with implicit patterns. In this case, the `ping` function simply outputs “Ping finished” to standard out and terminates if the first argument (the number of times to ping) is 0. If `N` is not 0, it sends a message to the `Pong` process, and then awaits a `pong` message back. The sent message is a tuple containing the atom `ping`, to indicate that this is a ping message, and a reference to the current process, obtained with the built-in `self` function. Once a ping has been sent and a pong received, it prints “Pong received”, and then recurses with one fewer ping left to send.

Now, let’s write `pong`:

```

1 pong() ->
2     receive
3         {ping, Ping} ->
4             Ping ! pong,
5             io:format("Ping received~n", [])
6     end,
7     pong().

```

Where `ping` starts with a send, `pong` instead starts with a receive. Once `pong` has received a message matching the pattern `{ping, Ping}` (remember, `Ping` is a variable because it starts with a capital letter), it sends a `pong` message back (`pong` is an atom in this case, not the function), and then prints “Ping received”. We’ve intentionally written this with the print after the send, to demonstrate concurrency.

If the `start` function we wrote above is run, one possible output is:

```

Ping received
Pong received
Ping received
Pong received
Ping received
Pong received
Ping received
Pong received
Ping received
Pong received
Ping received
Pong received
Ping finished

```

However, because the `pong` process sends its `pong` *before* printing that the `ping` was received, other orders are possible, such as this one:

```

Pong received
Ping received
Ping received
Pong received
Ping received
Pong received
Pong received
Ping received
Pong received
Ping received
Ping finished
Ping received

```

**Exercise 2.** What orders are *not* possible?

In this example, the pong process never actually terminates: only ping knew how many times to ping, so pong is left waiting endlessly for another ping that will never arrive. For cleanliness, we could instead add a terminate message like so:

```

1 ping(0, Pong) ->
2   io:format("Ping finished~n", []),
3   Pong ! terminate;
4
5 ping(N, Pong) ->
6   Pong ! {ping, self()},
7   receive
8     pong -> io:format("Pong received~n", [])
9   end,
10  ping(N - 1, Pong).
11
12 pong() ->
13  receive
14    terminate ->
15      io:format("Pong finished~n", []);
16
17    {ping, Ping} ->
18      Ping ! pong,
19      io:format("Ping received~n", []),
20      pong()
21  end.

```

Since pong does not recurse if terminate is received, it instead simply ends, terminating the process.

## 4 Processes as References

Erlang does not have mutable variables. But, surprisingly, they can be built with nothing but processes!

When representing mutable data in functional languages, we needed a way to put that data aside, separate from the program, in the heap ( $\Sigma$ ). But, concurrent processes are already “aside” and separate from one another, so all we actually need is a way for a process to store a piece of data like a single mapping in the heap, similarly to Haskell’s monads.

To achieve this, we will make a module which exports three functions: `ref/1`, `get/1`, and `put/2`. The `ref` function will generate a reference, like OCaml’s `ref`. The `get` function will retrieve the value stored in a reference, like OCaml’s `!`. The `put` function will store a value in a reference, like OCaml’s `:=`. The actual value used for the reference will be a process reference, to a process carefully designed to work this way.

The complete solution follows:

```

1 -module(refs).
2 -export([ref/1, refproc/1, get/1, put/2]).
3
4 ref(V) ->
5   spawn(refs, refproc, [V]).
6
7 refproc(V) ->
8   receive
9     {get, Return} ->
10      Return ! {refval, V},
11      refproc(V);
12
13     {put, Return, NV} ->
14      Return ! {refval, NV},
15      refproc(NV)
16   end.
17
18 get(Ref) ->
19   Ref ! {get, self()},
20   receive
21     {refval, V} -> V
22   end.
23

```

```

24 put(Ref, V) ->
25   Ref ! {put, self(), V},
26   receive
27     {refval, _} -> Ref
28   end.

```

The `ref` function spawns a new process, returning the process reference. The new process represents the reference, and runs the function `refproc`, with the initial value as its arguments. The `get` function sends a `get` message to the given process (which must be a reference process for this to work), and expects a `refval` message in response with the value stored in the reference. The `put` function sends a `put` message to the given process, containing the value to put in the reference, and also waits for a `refval` message. `put` doesn't actually care about the value returned by `refval`; it's only used to make sure that the message has been received and acted on before the current process continues.

The `refproc` function contains all of the interesting behavior, as it is the function used by the actual reference process. `refproc` must be exported because of how `spawn` works—there are ways to get around “polluting” the exported names in this way, but they're not important for our purposes. `refproc`'s behavior is quite similar regardless of whether it receives a `get` message or a `put` message: it returns a value and then recursively calls `refproc` again. The difference is in which value. The value stored in the reference is in the (immutable) `V` variable. With `get`, it returns that value, and then recurses with the same value. With `put`, it instead expects a new value, `NV`, and returns and recurses with `NV` instead of `V`. In this way, although no variables are mutable, the reference itself is, since if it receives a `put` message, then it will respond to future `get` messages with the new value, until another `put` is received.

In shared-memory concurrency, two tasks may access the same mutable memory, and each may mutate it. With these references, we have in fact implemented shared-memory concurrency on top of message-passing concurrency: if two processes each have a reference to the reference process, then either may mutate its value, and both can see the other's mutations. The fact that it is then extremely difficult to guarantee that the processes mutate things in the correct order is the usual argument for using message-passing concurrency instead of shared-memory concurrency in the first place, but this form of references demonstrates that neither is more powerful than the other.

## 5 Processes as Objects

If you're accustomed to object-oriented programming, you've probably noticed that references from the reference module above act a lot like objects with a `get` and `put` method. Indeed, we can extend this metaphor to implement objects with only processes with immutable variables. For instance, this module implements a reverse polish notation calculator object very similar to the one we wrote for the Smalltalk segment of Module 1:

```

1 -module(rpncalc).
2 -export([newrpn/0, rpn/1, push/2, binary/2, add/1, sub/1, mul/1, divide/1]).
3
4 newrpn() ->
5   spawn(rpncalc, rpn, [[]]).
6
7 rpn(Stack) ->
8   receive
9     {push, Return, V} ->
10      Return ! {rpnval, V},
11      rpn([V | Stack]);
12
13     {op, Return, F} ->
14      rpnop(Stack, F, Return)
15   end.
16
17 rpnop([R, L | Rest], F, Return) ->
18   V = F(L, R),
19   Return ! {rpnval, V},
20   rpn([V | Rest]).
21
22 push(RPN, V) ->
23   RPN ! {push, self(), V},
24   receive

```

```

25     {rpnval, _} -> V
26     end.
27
28 binary(RPN, F) ->
29     RPN ! {op, self(), F},
30     receive
31     {rpnval, V} -> V
32     end.
33
34 add(RPN) ->
35     binary(RPN, fun(L, R) -> L + R end).
36
37 sub(RPN) ->
38     binary(RPN, fun(L, R) -> L - R end).
39
40 mul(RPN) ->
41     binary(RPN, fun(L, R) -> L * R end).
42
43 divide(RPN) ->
44     binary(RPN, fun(L, R) -> L / R end).

```

An RPN’s sole field, the stack, is represented by the `Stack` variable of the `rpn` function, which is the function run by an RPN process. The `push` and `binary` functions send an RPN process messages corresponding to one of its two supported “methods”: `push` and `op`. The `op` message carries a function, representing the binary operation to perform, so like in the Smalltalk version, specific operations can be implemented in terms of it.

Again, there are only immutable variables and processes, but the fact that sending and receiving messages establishes a sequence allows us to emulate more sophisticated features. Erlang has several libraries implementing more elegant object orientation, but still using processes as objects. Erlang is designed for programs to have thousands of processes, so it’s common to mix these styles as well; for instance, fields can be implemented as references which are in turn implemented as processes as in the previous section.

## 6 Implementation

Operating systems implement processes and threads: in OS terms, two processes do not share memory, but two threads within the same process do share memory. Erlang uses the term “process” because Erlang’s processes do not share memory. However, using operating system processes to implement Erlang processes would result in catastrophically poor performance! Indeed, even using operating system threads to implement Erlang processes would be similarly fraught. The reason is simply that switching between threads or processes in an operating system—that is, context switching—is an expensive process.

Instead, highly-concurrent languages such as Erlang use so-called *green threads*. Basically, the Erlang interpreter must implement its own form of thread switching, and maintain the stack for each thread as a native data structure in the host language. When an Erlang process executes a `receive` and it does not immediately have a matching message available to act upon, Erlang instead sets aside the thread for that process and loads a stack for another process; in effect, it performs the same kind of context switch that an operating system performs, but with much less context to switch. It runs as many operating system threads as there are CPU cores available, but each one can switch between many green threads. It is not uncommon for an Erlang process to have tens of thousands of processes, so keeping green threads light is extremely important.

To know which processes are available to run, an Erlang process must also be able to pattern match very quickly. Typically, a process that is waiting for a message has its pattern stored, and when another process *sends* it a message, it performs a pattern match immediately, to determine if the other process can be awoken.

The other major implementation roadblock to message-passing concurrency is the actual message passing. In languages with mutable values, it is necessary to copy the message being passed, so that two processes cannot see the same mutable memory (which would be shared-memory concurrency, not message-passing concurrency). Erlang largely sidesteps this issue by being fundamentally immutable, so that it’s harmless to pass around values in any form. Two processes can share a pointer to a value if neither will actually mutate that value.

## 7 Miscellany

As mentioned previously, we won't mention how to actually solve any interesting problems with concurrency; instead, you've hopefully realized, particularly through the examples of processes as references and objects, that like the  $\lambda$ -calculus before it,  $\pi$ -calculus and Erlang demonstrate that a surprising amount of power is unleashed merely by having concurrency, before we've even reached actual algorithms. Most programming language semantics that aren't designed specifically to address concurrency simply ignore it, because adding non-determinism affects all other aspects of a semantics.

Erlang does actually support some (very limited) mutable data structures, but they may not be sent in a message.

The message-reply style we used in all of our examples was fairly brittle, in that we used a specific atom as the expected reply, but there's nothing to stop another process from sending the same atom. Erlang supports creating unique values, called "references" to create needless confusion, with the built-in `create_ref` function. Usually, two processes which communicate would exchange such unique references to make sure that they're receiving messages from the process they thought they were speaking to.

Most Erlang programs are written in a so-called "let it crash" style. That is, instead of trying to anticipate all forms of errors, code is written to simply re-spawn processes that fail unexpectedly. Since processes are mostly independent of each other, large systems can operate even with major bugs. In the Erlang shell, you can re-compile modules, and swap out processes using the old module for processes using the new version, and it is thus often possible to fix bugs in a running system with no downtime. Many Erlang proponents cite this style as the major advantage of Erlang.

## 8 Fin

In the next (and final) module, we will look at how our mathematical model of programming languages interfaces with the real world, through models of systems programming.

## References

- [1] Robin Milner, Joachim Parrow, and David Walker. A calculus of mobile processes, I. *Information and computation*, 100(1):1–40, 1992.

### Rights

Copyright © 2020–2023 Gregor Richards.  
This module is intended for CS442 at University of Waterloo.  
Any other use requires permission from the above named copyright holder(s).