

INTRODUCTORY GUIDE TO MESSAGE PASSING IN DISTRIBUTED COMPUTER SYSTEMS

Florian Willich
Hochschule für Technik und Wirtschaft Berlin
University of Applied Sciences Berlin
Course: Distributed Systems
Lecturer: Prof. Dr. Christin Schmidt

June 7, 2015

Abstract

Message passing in distributed systems is a model to exchange messages within a process pair by making use of several standards and implementation details. Those have been developed to offer the right message passing models for the different areas of applications. The programming language Erlang natively supports an asynchronous message passing model which makes the implementation of concurrent applications transparent to the software developer.

Contents

1	MESSAGE ORIENTED COMMUNICATION	3
1.1	INTRODUCTION	3
1.2	BASIC REQUIREMENTS IN THEORY	3
1.3	REQUIREMENTS PROVIDER IN PRACTICE	4
1.4	ASYNCHRONOUS VS. SYNCHRONOUS MESSAGE PASSING	4
2	MESSAGE-PASSING INTERFACE STANDARD	5
3	MESSAGE PASSING IN ERLANG	6
3.1	Introduction	6
3.2	Concurrency in Erlang	6
3.3	Implementation of a Simple Diffie-Hellman Key Exchange in Erlang	7
4	DEFINITION OF MESSAGE PASSING IN DISTRIBUTED SYSTEMS	10
A	my_math Erlang Module	11
B	Acronyms	12

1 MESSAGE ORIENTED COMMUNICATION

1.1 INTRODUCTION

Communication in distributed systems defines a distributed system itself: N systems execute one or more tasks by computing and communicating with each other. Every time systems shall communicate with each other to provide services, a proficient communication model has to be implemented. There are several models of communication in distributed systems such as Remote Procedure Calls (RPC) or object method invocation (Tanenbaum & Steen, 2007, chap. 4.3 / p. 203). With this paper I'll introduce you to the message passing model in distributed systems.

As human beings we already have a deep understanding of different models of message passing. When you read this words, I am passing a message to you, which seems to be a trivial thing to people who can read. In more philosophical terms, now that you read those words, a message is passed from one entity (me, the writer) to another entity (you, the reader) which is obviously non trivial considering all the assumptions we would've to make to actually perform this message passing model of human beings.

Message passing in distributed systems e.g. computer science is based on messages, composed of bit strings, exchanged within a process pair - which would be the equivalent to the entity pair. It is important to understand that message passing is a model designed for Inter-Process Communication (IPC). Where those processes are located, on one or on different systems, is subsequent to the provided functionality.

1.2 BASIC REQUIREMENTS IN THEORY

The following items are the basic requirements (in theory) for a system to provide the functionality of message passing:

- **Connectivity:** A connection to communicate has to be established between the process pair
- **Ability:** Each process has to be able to receive or send messages
- **Integrity:** Send messages have to be delivered as is
- **Intelligibility:** The receiving process has to be able to interpret the message as intended
- **Executability:** The delivered message has to lead to an execution of the desired instructions

1.3 REQUIREMENTS PROVIDER IN PRACTICE

To provide the above mentioned requirements there have been developed several message passing models resulting in well defined standards and concrete implementations in the last decades. To make it easier to understand the relations between the theoretical requirements and a concrete implementation, I will take the Transmission Control Protocol / Internet Protocol (TCP/IP) which uses several socket primitives as example (Tanenbaum & Steen, 2007, chap 4.3.1 on p. 141 - 142):

- **Connectivity:** A communication end point has to be provided so the application can write data to and read data from: TCP/IP provides the *Socket* primitive to create that end point and also the *Bind* primitive to bind a local address to that socket. The *Connect* primitive then provides the functionality to establish a connection.
- **Ability:** TCP/IP provides the two primitives *Send* and *Receive* to simply send and receive data over the connection. (A little but important detail: The Operating System (OS) has to reserve local memory to provide a buffer for the in- and outgoing messages)
- **Integrity:** TCP/IP provides several mechanisms to ensure that there has been no data loss when sending or receiving messages. On the other hand this makes the protocol more slow than other protocols such as the Real-Time Transport Protocol (RTP) or the User Datagram Protocol (UDP).

The requirements of **Intelligibility** and **Executability** do not depend on TCP/IP and thus other standards have to take place to provide those requirements such as Extensible Markup Language (XML) or JavaScript Object Notation (JSON) to structure the messages in a standardized manner and last but not least the message parsing implementation to call the desired instructions.

There are other primitives than the above mentioned to provide the functionality of TCP/IP. If you're interested in reading more, I recommend you the referenced book (Tanenbaum & Steen, 2007, ch. 4.3.1 on p. 141).

It is also important to understand, that this was just a round-up of how you could basically make use of the socket primitives to implement message passing in your application. There is plenty to discuss about how message passing is realized such as the Open Systems Interconnection Reference Model (OSI) which was developed by the International Organization for Standardization (ISO) to model the different layers of network oriented communication (Tanenbaum & Steen, 2007, ch. 4.1.1 on p. 116).

1.4 ASYNCHRONOUS VS. SYNCHRONOUS MESSAGE PASSING

Now that we discussed a protocol to handle the socket primitives one major difference in message passing models is the whether we use asynchronous or

synchronous message passing. Passing a message synchronously means that the called message send routine returns after the request has been successfully transmitted (or an error occurred). Receiving a message synchronously means that the called message receiving routine reads a specific amount of bytes from the socket and returns that message.

To pass or receive our message asynchronously another application layer or middleware has to take place. The call of a send routine can either return when the middleware (Tanenbaum & Steen, 2007, ch. 4.1 on p. 125):

1. took over the transmission of the request.
2. successfully send the request to the receiver.
3. successfully send the request to the receiver, assuring this by a corresponding message.

The middleware can also perform some preparatory work e.g. separating stand alone request strings and parsing them into data structure and returning this more easy to parse data when calling the receive routine of the message passing middleware.

One key model to provide such middleware functionalities are message queues. A message queue is a queue to store messages by the First In - First Out (FIFO) principle. This enables the application to asynchronously send and receive messages by offering an incoming message queue and an outgoing message queue. This makes all the send and receive mechanisms fully transparent to the application (Tanenbaum & Steen, 2007, ch. 4.3.2 on p. 145 - 147)

2 MESSAGE-PASSING INTERFACE STANDARD

The Message-Passing Interface (MPI) is a message-passing library interface specification, made for high performance and scale able distributed systems where high-efficiency is needed. The MPI standard was designed by the Message-Passing Interface Forum (MPIF) which is an open group with representatives from many organizations. The current version of the MPI standard is MPI-3.0 (Message Passing Interface Forum, 2012, ch. Abstract/ii & Acknowledgements/xx & 1.1 on p. 1).

The MPIF aims to offer a standard which establishes a practical, portable, efficient and flexible way of implementing message-passing in various high-level programming languages (e.g. C, C++, Fortran) (Message Passing Interface Forum, 2012, ch. History/iii). Furthermore, the MPI simplifies the communication primitives and brings them to an abstraction level to perfectly fit the programmers needs of writing efficient and clean code for such distributed systems

(Tanenbaum & Steen, 2007).

TCP does not fit those requirements. While the socket primitives *read* and *write* are sufficient for several general-purpose protocols, (managing the communication across networks) they are insufficient for high-speed interconnection networks such as super computers or server clusters. The MPI standard offers a set of functions and datatypes with which the software developer is able to explicitly execute synchronous and asynchronous message passing routines (Tanenbaum & Steen, 2007, ch. 4.3 on p. 143).

3 MESSAGE PASSING IN ERLANG

3.1 Introduction

Erlang is a functional, declarative programming language that was written for the need of real-time, non-stop, concurrent, very large and distributed system applications (Armstrong, 1996, chap. 1 / p. 1). While the language was originally designed by Joe Armstrong, Robert Virding and Mike Williams for Ericsson (a Swedish telecommunications provider) in 1986, it finally became open source in 1998, thanks to the open source initiatives lead by Linux (Däcker, 2000, chap. 8 on p. 39).

Since there is a lot confusion with the naming of Erlang: Joe Armstrong, who started to design Erlang by adding functionality to Prolog, named it after the Danish mathematician Agner Krarup Erlang (creator of the Erlang loss formula) following the tradition of naming programming languages after dead mathematicians (Däcker, 2000, chap. 4.1 on p. 13).

Erlang uses a native, asynchronous message passing model to communicate between light weight processes, also called actors (Armstrong, 1996, chap. 1 on p. 1). Erlang does not make heavy use of your OS to provide the described concurrency model which implicitly means that Erlang decouples the underlying OS and thus providing a cross-platform and transparent message passing model (Armstrong, 1996, chap. 1 & 3 on p. 1 - 3).

3.2 Concurrency in Erlang

Erlang provides semantics and built-in functions to parallelize your applications by message passing (Ericsson AB, 2015, ch. 4.3 on p. 95 - 104):

- **spawn**(*Module*, *Exported Function*, *List of Arguments*)
Function which creates a new actor, by running the *Exported Function* with the *List of Arguments* located in the *Module* (a set of functions located in one file) and returns the Process Identifier (pid), which uniquely identifies

an actor, of it.

- The **receive** construct that allows the function that is executed by an actor to receive messages by using a message queue.
- The **!** operator which sends the right handed term to the left handed pid. The right handed term is the message we send.
- **self()**
Function that returns the pid of the actor who executes the function.

Although the above described semantics and functions require a little knowledge on the Erlang programming language it was important for me to show you that Erlang is not only offering you an easy to understand interface for concurrency, it rather is a concurrent functional programming language. I recommend you to have a look into the referenced official Erlang documentation or to *learn you some Erlang* on www.learnyousomeerlang.com.

3.3 Implementation of a Simple Diffie-Hellman Key Exchange in Erlang

```
1 %%% @author      Florian Willich
2 %%% @copyright   The MIT License (MIT) for more information see
3 %%%              http://opensource.org/licenses/MIT
4 %%% @doc         This module represents a simple implementation of
5 %%%              the
6 %%%              Diffie-Hellman key exchange algorithm.
7 %%% @end
8 %%% Created 2015-06-06
9
10 -module(diffie_hellman).
11 -author("Florian Willich").
12 -compile(export_all).
13
14 %%% @doc Public data represents the data which is publicly shared
15 %%%      within two communication partners when exchanging key by
16 %%%      the
17 %%%      Diffie-Hellman key exchange algorithm.
18 %%%      p: The public prime number
19 %%%      g: The public prime number (1 ... p - 1)
20 %%%      componentKey: The computed component key
21 %%%      pid: The pid of the one who instantiated this record
22 -record(publicData, {p, g, componentKey, pid}).
23
24 %%% @doc Returns the value G to the power of MySecretKey Modulo P
25 %%%      which is the public component key for the Diffie-Hellman
26 %%%      key
27 %%%      exchange algorithm.
28 %%%      For more information see http://goo.gl/pzdiH
29 %%% @end
30
31 -spec computeMyPublicComponentKey(pos_integer(), pos_integer(),
32                                   pos_integer()) -> pos_integer().
```

```

28 computeMyPublicComponentKey(P, G, MySecretKey) ->
29   my_math:pow(G, MySecretKey) rem P.
30
31 %%% @doc Returns the value ComponentKey to the power of
    MySecretKey
32 %%%      Modulo P which is the private shared key for the
33 %%%      Diffie-Hellman key exchange algorithm.
34 %%%      For more information see http://goo.gl/pzdiH.
35 %%% @end
36 -spec computeSharedPrivateKey(pos_integer(), pos_integer(),
    pos_integer()) -> pos_integer().
37 computeSharedPrivateKey(P, ComponentKey, MySecretKey) ->
38   my_math:pow(ComponentKey, MySecretKey) rem P.
39
40 %%% @doc Starts the Diffie-Hellman key exchange algorithm by
    taking P
41 %%%      (a prime number), G (1 ... P - 1), MySecretKey is the
    secret
42 %%%      integer of the one who executes this function and the
43 %%%      PartnerPID which is the PID of the communication partner
44 %%%      with you a key exchange shall be initiated. This function
45 %%%      sends the term {startKeyExchange, PublicData} to the
46 %%%      PartnerPID where PublicData is of type publicData.
47 %%%      Afterwards, the function starts a receive construct which
    is
48 %%%      receiving the following:
49 %%%      {componentKey, PublicData}:
50 %%%      The message including all information needed for
51 %%%      computing the private shared key and then prints it
    out.
52 %%%      UnexpectedMessage:
53 %%%      Prints out any unexpected incoming message and calls
    a
54 %%%      recursion.
55 %%%      After 3000 milliseconds:
56 %%%      The function will return timeout.
57 %%% @end
58 -spec startKeyExchange(pos_integer(), pos_integer(), pos_integer(),
    term()) -> term().
59 startKeyExchange(P, G, MySecretKey, PartnerPID) ->
60   MyComponentKey = computeMyPublicComponentKey(P, G, MySecretKey),
61   MyPublicData = #publicData{p = P, g = G, componentKey =
    MyComponentKey, pid = self()},
62   PartnerPID ! {startKeyExchange, MyPublicData},
63
64   receive
65
66     {componentKey, #publicData{p = P, g = G, componentKey =
    PartnerComponentKey, pid = PartnerPID}} ->
67       PrivateSharedKey = computeSharedPrivateKey(P,
    PartnerComponentKey, MySecretKey),
68       printSharedPrivateKey(self(), PrivateSharedKey);
69
70   UnexpectedMessage ->
71     printUnexpectedMessage(UnexpectedMessage),
72     startKeyExchange(P, G, MySecretKey, PartnerPID)
73

```



```

74  after 3000 ->
75      timeout
76
77  end.
78
79  %%% @doc Listens on Messages to start the Diffie-Hellman key
      exchange
80  %%%      by starting the following receive construct:
81  %%%      {startKeyExchange, PublicData}:
82  %%%          The message including all information needed to start
83  %%%          the key exchange by computing the own public data
      which
84  %%%          will then be send to the PartnerPID as follows:
85  %%%          {componentKey, MyPublicData}.
86  %%%          Afterwards, the private shared key will be printed
      out
87  %%%          and the function calls a recursion.
88  %%%      terminante:
89  %%%          Prints out that this function terminates with the
90  %%%          executing PID and returns ok.
91  %%%      UnexpectedMessage:
92  %%%          Prints out any unexpected incomping message and calls
      a
93  %%%          recursion.
94  %%% @end
95  -spec listenKeyExchange(pos_integer()) -> term().
96  listenKeyExchange(MySecretKey) ->
97      receive
98
99      {startKeyExchange, #publicData{p = P, g = G, componentKey =
      PartnerComponentKey, pid = PartnerPID}} ->
100      MyComponentKey = computeMyPublicComponentKey(P, G,
      MySecretKey),
101      MyPublicData = #publicData{p = P, g = G, componentKey =
      MyComponentKey, pid = self()},
102      PartnerPID ! {componentKey, MyPublicData},
103      PrivateSharedKey = computeSharedPrivateKey(P,
      PartnerComponentKey, MySecretKey),
104      printSharedPrivateKey(self(), PrivateSharedKey),
105      listenKeyExchange(MySecretKey);
106
107      terminate ->
108      io:format("~p terminates!~n", [self()]),
109      ok;
110
111      UnexpectedMessage ->
112      printUnexpectedMessage(UnexpectedMessage),
113      listenKeyExchange(MySecretKey)
114
115  end.
116
117  %%% @doc Prints out the UnexpectedMessage as follows:
118  %%%      Received an unexpected message: 'Unexpected Message'
119  %%% @end
120  printUnexpectedMessage(UnexpectedMessage) ->
121  io:format("Received an unexpected message: ~p~n", [
      UnexpectedMessage]).

```

```

122
123 %%% @doc Prints out the shared private key as follows:
124 %%%      'PID': The shared private Key is: 'SharedKey'
125 %%% @end
126 printSharedPrivateKey(PID, SharedKey) ->
127   io:format("~p: The shared private Key is: ~p~n", [PID, SharedKey]
128             ).
129 %%% @doc Starts a key exchange example by spawning the Alice
130 %%%      process,
131 %%%      which executes the listenKeyExchange function with
132 %%%      MySecretKey = 15, and the Bob process, which executes the
133 %%%      startKeyExchange function with P = 23, G = 5,
134 %%%      MySecretKey = 6 and PartnerPID = Alice. Returns Alice and
135 %%%      Bob.
136 startExample() ->
137   Alice = spawn(diffie_hellman, listenKeyExchange, [15]),
138   Bob = spawn(diffie_hellman, startKeyExchange, [23, 5, 6, Alice]),
139   {Alice, Bob}.
140
141 %%% @doc Starts a key exchange remote example by spawning the
142 %%%      Alice
143 %%%      process, which executes the listenKeyExchange function
144 %%%      with
145 %%%      MySecretKey = 15, and the Bob process, located on the
146 %%%      RemoteNode, which executes the startKeyExchange function
147 %%%      with P = 23, G = 5, MySecretKey = 6 and PartnerPID =
148 %%%      Alice.
149 %%%      Returns Alice and Bob.
150 %%% @end
151 startRemoteExample(RemoteNode) ->
152   Alice = spawn(RemoteNode, diffie_hellman, listenKeyExchange, [
153               15]),
154   Bob = spawn(diffie_hellman, startKeyExchange, [23, 5, 6, Alice]),
155   {Alice, Bob}.

```

4 DEFINITION OF MESSAGE PASSING IN DISTRIBUTED SYSTEMS

Message passing in distributed systems is a model to exchange messages within a process pair by making use of several standards and implementation details. The specifically used message passing model can diverge extremely in its provided functionality and defines how to provide the connection and send or receive messages whereas the physical conditions and the implementation of executing the desired instructions is undefined.

A my_math Erlang Module

```
1 %%% @author      Florian Willich
2 %%% @copyright   The MIT License (MIT) for more information see
3 %%%              http://opensource.org/licenses/MIT
4 %%% @doc         This is my math module for mathematical functions
                    not provided
5 %%%              by the erlang standard library.
6 %%% @end
7 %%% Created 2015-06-06
8
9 -module(my_math).
10 -author("Florian Willich").
11 -export([pow/2]).
12
13 %%% Returns the value of Base to the power of Exponent.
14 %%% If Base and Exponent is 0 the function returns
    undefinedArithmeticExpression.
15 %%% The motivation to implement this function was that there is no
    erlang
16 %%% standard library pow function returning an integer.
17 -spec pow(integer(), integer()) -> number().
18 pow(0, 0) ->
19     undefinedArithmeticExpression;
20
21 pow(Base, 0) ->
22     case Base < 0 of
23         true -> -1;
24         false -> 1
25     end;
26
27 pow(Base, Exponent) ->
28     case Exponent < 0 of
29         true -> 1 / pow(Base, -Exponent, 0);
30         false -> pow(Base, Exponent, 0)
31     end.
32
33 %%% Returns the value of Base to the power of Exponent. Acc should
    be 0 for
34 %%% initiating computation.
35 %%% If Base and Exponent is 0 the function returns
    undefinedArithmeticExpression.
36 %%% The motivation to implement this function was that there is no
    erlang
37 %%% standard library pow function returning an integer.
38 -spec pow(pos_integer(), non_neg_integer(), non_neg_integer()) ->
    integer().
39 pow(0, 0, _) ->
40     undefinedArithmeticExpression;
41
42 pow(_, 0, Acc) -> Acc;
43
44 pow(Base, Exponent, 0) ->
45     pow(Base, Exponent - 1, Base);
46
47 pow(Base, Exponent, Acc) ->
48     pow(Base, Exponent - 1, Acc * Base).
```

B Acronyms

Acronyms

FIFO First In - First Out. 5

IPC Inter-Process Communication. 3

ISO International Organization for Standardization. 4

JSON JavaScript Object Notation. 4

MPI Message-Passing Interface. 5

MPiF Message-Passing Interface Forum. 5

OS Operating System. 4, 6

OSI Open Systems Interconnection Reference Model. 4

pid Process Identifier. 6, 7

RPC Remote Procedure Calls. 3

RTP Real-Time Transport Protocol. 4

TCP/IP Transmission Control Protocol / Internet Protocol. 3, 4

UDP User Datagram Protocol. 4

XML Extensible Markup Language. 4

References

- Armstrong, Joe. 1996. *Erlang - A survey of the language and its industrial applications*. In: *Proceedings of the symposium on industrial applications of Prolog (INAP96)*. <http://www.erlang.se/publications/inap96.ps>. The 9'th Exhibitions and Symposium on Industrial Applications of Prolog. 16-18, October 1996. Hino, Tokyo Japan. [Online. Accessed 5th May 2015].
- Däcker, Bjarne. 2000. *Concurrent functional programming for telecommunications: A case study of technology introduction*. <http://www.erlang.se/publications/bjarnelic.ps>. Licenciate Thesis, Department of Teleinformatics. TRITA-IT AVH 00:08, ISSN 1403-5286. Royal Institute of Technology, Stockholm, Sweden. [Online. Accessed 20th May 2015].
- Ericsson AB. 2015. *Erlang/OTP System Documentation 6.4*. <http://www.erlang.org/doc/pdf/otp-system-documentation.pdf>. [Online. Accessed 31th May 2015].
- Hebert, Fred. 2013. *Learn You Some Erlang for Great Good! : A Beginner's Guide*. <http://learnyousomeerlang.com/content>. No Starch Press. [Online. Accessed 20th May 2015].
- Message Passing Interface Forum. 2012. *MPI: A Message-Passing Interface Standard Version 3.0*. <http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf>. [Online. Accessed 20th May 2015].
- Tanenbaum, Andrew S., & Steen, Maarten Van. 2007. *Distributed Systems: Principles and Paradigms (Second Edition)*. Pearson Prentice Hall. ISBN 0-13-239227-5.