# An Erlang High Performance TCP/IP Stack

Javier París    Víctor M. Gulías    Alberto Valderruten

<{paris,gulias,valderruten}@dc.fi.udc.es>

# Índice

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction
Design
Optimization
Issues
Tests
Distribution

# Motivation

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

**Introduction**

Design

Optimization
Issues

Tests

Distribution

- Problem:
    - Streaming video server ir Erlang.
    - Uses protocols like HTTP and RTSP => TCP.
    - Connections usually last for hours.
    - Cut in the connection => playing stops.
    - Clients are usually closed settop boxes.
- Tcp must be fault tolerant in the server side.

# Índice

**An Erlang High Performance TCP/IP Stack**

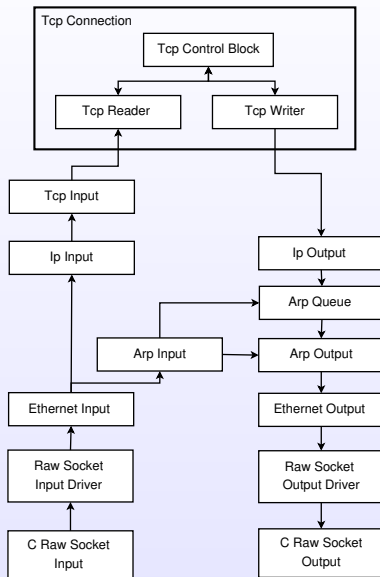**Javier París, Víctor M. Gulías, Alberto Valderruten**

Introduction

**Design**

Optimization Issues

Tests

Distribution

# Overall Design

Proccess Structure:

An Erlang High Performance TCP/IP Stack

Javier París, Víctor M. Gulías, Alberto Valderruten

Introduction
Design
Optimization Issues
Tests
Distribution

# Network Interface

An Erlang High Performance TCP/IP Stack

Javier París, Víctor M. Gulías, Alberto Valderruten

Introduction

**Design**

Optimization Issues

Tests

Distribution

- Erlang does not support file operation on devices => Linked-in driver.
- Linked-in to avoid unnecessary data copies.
- Interface with the card is a Linux Raw Socket.
- Network Card must be in promiscous mode.
- Data is read and written directly to the interface.

# Server Model

**An Erlang High Performance TCP/IP Stack**

**Javier París, Víctor M. Gulías, Alberto Valderruten**

**Introduction**

**Design**

**Optimization Issues**

**Tests**

**Distribution**

- Two proccesses for each stateless protocol:

```
init_writer(Conf) ->
  register(protocol_writer, self()),
  writer_loop(Conf).

writer_loop(Conf) ->
  receive
    {send, Packet, Dst} ->
      {ok, Bin_Packet}=build_packet(Packet,
                                     Dst),
      send_packet(Bin_Packet, Dst, Conf);
  end,
  writer_loop(Conf).
```

# Server Model

An Erlang High Performance TCP/IP Stack

Javier París, Víctor M. Gulías, Alberto Valderruten

Introduction

Design

Optimization Issues

Tests

Distribution

- The reader is similar:

```erlang
init_reader(Conf) ->
  register(protocol_reader, self()),
  reader_loop(Conf).


reader_loop(Conf) ->
  receive
    {recv, Packet} ->
      case catch decode(Packet, Conf) of
        {ok, Upper_protocol, Data} ->
          Upper_protocol:recv(Data);
        {error, Error} ->
          {error, Error};
      end
  end,
  reader_loop(Conf).
```

# Tcp Design

An Erlang High Performance TCP/IP Stack

Javier París, Víctor M. Gulías, Alberto Valderruten

Introduction

Design

Optimization Issues

Tests

Distribution

Tcp is stateful. Must keep state for each connection:

- Reader process.
- Writer process.
- TCB(Tcp Control Block): state and synchrnization.

# State Behaviour

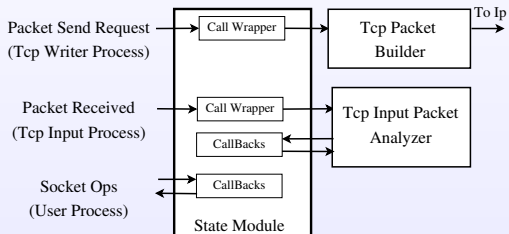State beaviour is implemented using callbacks and wrappers:

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction

**Design**

Optimization
Issues

Tests

Distribution

# Índice

**An Erlang High Performance TCP/IP Stack**

Javier París,
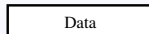Víctor M. Gulías,
Alberto Valderruten

Introduction

Design

**Optimization Issues**

Tests

Distribution

# Scatter/Gather

**An Erlang High Performance TCP/IP Stack**

**Javier París, Víctor M. Gulías, Alberto Valderruten**

**Introduction**

**Design**

**Optimization Issues**
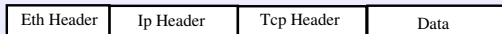
**Tests**

**Distribution**

Functional Scatter Gather:

| Data |
|------|

<<Data>>

| Tcp Header | Data |
|------------|------|

[<<Tcp Header>>,<<Data>>]

| Ip Header | Tcp Header | Data |
|-----------|------------|------|

[<<Ip Header>>, <<Tcp Header>>,<<Data>>]

| Eth Header | Ip Header | Tcp Header | Data |
|------------|-----------|------------|------|

[<<Eth Header>>, <<Ip Header>>, <<Tcp Header>>,<<Data>>]

# Checksum

**An Erlang High
Performance
TCP/IP Stack**

**Javier París,
Víctor M. Gulías,
Alberto
Valderruten**

**Introduction**

**Design**

**Optimization
Issues**

**Tests**

**Distribution**

Checksum was done in several ways:

- Deconstructing a binary one byte at a time: 13MB/s

```
checksum(<<Num:16/integer, Remainder/binary>>
          Csum) ->
    checksum(Remainder, Csum + Num);
```

- Deconstructing a binary eight bytes at a time:
  28MB/s

```
checksum(Bin = <<N1:16/integer,N2:16/integer,
                 N3:16/integer,N4:16/integer,
                 N5:16/integer,N6:16/integer,
                 N7:16/integer,N8:16/integer,
                 Rem/binary>>, Csum)
                   when size(Bin) >= 16 ->
    checksum(Rem, Csum+N1+N2+N3+N4+N5+N6+
                      N7+N8);
```

# Checksum

**An Erlang High Performance TCP/IP Stack**

Javier París, Víctor M. Gulías, Alberto Valderruten

**Introduction**

**Design**

**Optimization Issues**

**Tests**

**Distribution**

- Using the same binary eight bytes at a time: 32MB/s
```
checksum(Bin, Position) when size(Bin) >= 8 ->
<<_:(Position)/integer, N1:16/integer,
        N2:16/integer, ..., _/binary>> = Bin,
        checksum(Rem, Csum+N1+N2..+N8);
```
- Using a C linked-in driver: 136MB/s

# Índice

**An Erlang High Performance TCP/IP Stack**

Javier París, Víctor M. Gulías, Alberto Valderruten
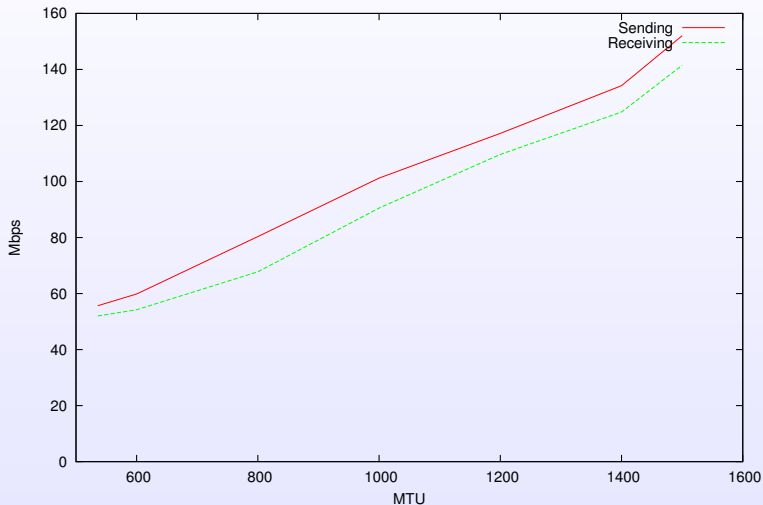
Introduction

Design

Optimization Issues

**Tests**

Distribution

# Throughput

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Throughput Transfering a 1GB File

# Concurrency

**An Erlang High Performance TCP/IP Stack**

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction

Design

Optimization
Issues

Tests

Distribution

Throughput under Concurrent Connections
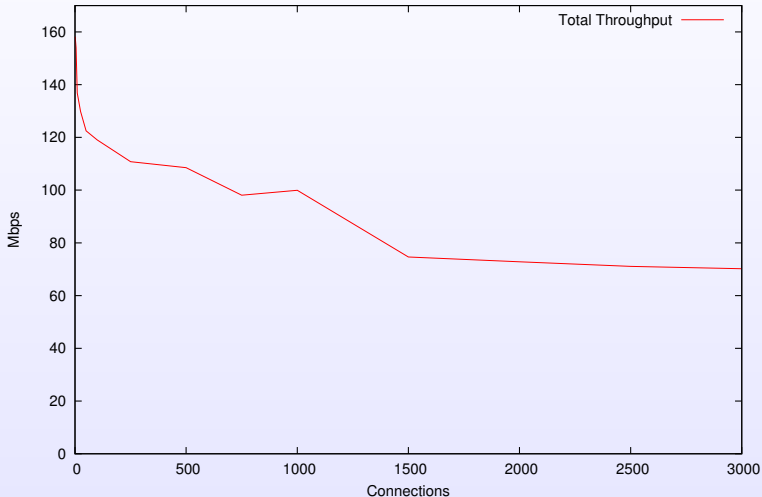
# Índice

**An Erlang High
Performance
TCP/IP Stack**

**Javier París,
Víctor M. Gulías,
Alberto
Valderruten**

Introduction
Design
Optimization
Issues
Tests
**Distribution**

# Distribution

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction
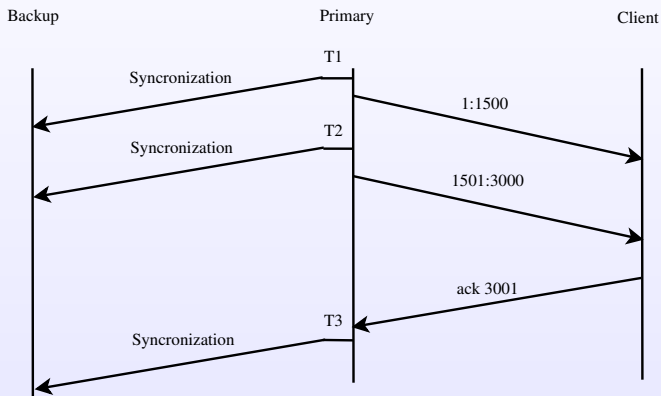
Design

Optimization
Issues

Tests

Distribution

- Aim: Make it possible to recover ongoing Tcp connections of failed machines in backup nodes.
- Erlang applications are often distributed => Try to take advantage of distributed information about the state to ease the cost.

# Distribution: Sending Data

An Erlang High Performance TCP/IP Stack

Javier París, Víctor M. Gulías, Alberto Valderruten

Introduction

Design

Optimization Issues

Tests

Distribution

First Approach: Make a synchronization before everey packet sent and after every packet received.

# Distribution: Optimizing Send

**An Erlang High Performance TCP/IP Stack**

**Javier París, Víctor M. Gulías, Alberto Valderruten**

**Introduction**

**Design**

**Optimization Issues**

**Tests**

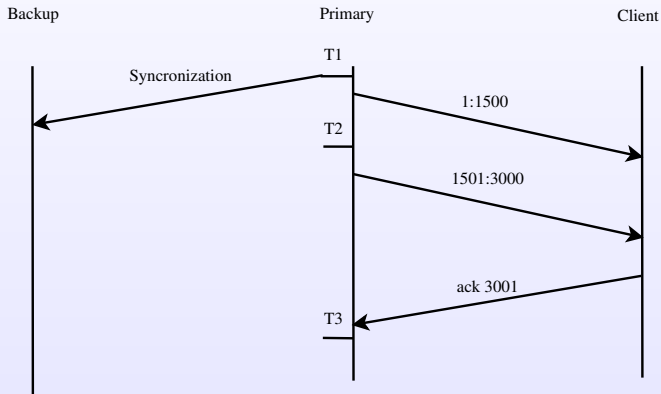**Distribution**

Optimization:

- Data that must be synchronized at all times: Sequence Numbers and Window Information.
- Client also has this information at all times.
- An out of window packet must be answered with an ack (RFC 793) => The backup can recover this information by sending an out of window packet to the client.
- The client can use initial sequence number to generate an invalid sequence number. When the initial sequence number is about to enter window again the primary node must warn the backups => Forced Synchronization only every 4GB.

# Distribution: Optimizing Send

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction
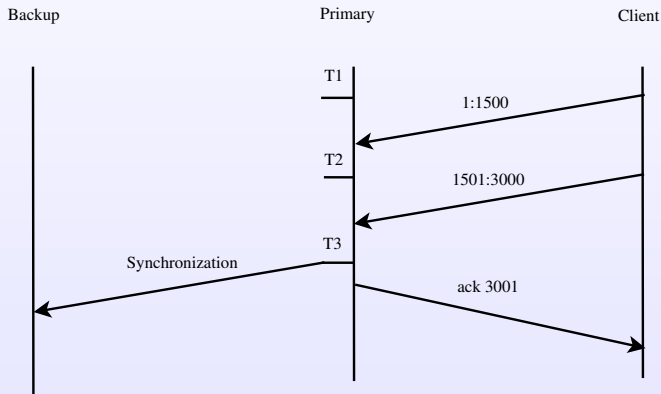
Design

Optimization
Issues

Tests

Distribution

- The application must continue where it left in the failed node. The stack may help by saying how many bytes have been transmitted.

# Distribution Reception

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction

Design

Optimization
Issues

Tests

Distribution

Harder to Optimize:

- Buffered received data cannot be recovered.
- Sincronization is needed every time an ack is received.

# Distribution

An Erlang High
Performance
TCP/IP Stack

Javier París,
Víctor M. Gulías,
Alberto
Valderruten

Introduction

Design

Optimization
Issues

Tests

Distribution

Future Challenges:

- Failure detection and collision avoidance between master and slaves.
- How to physically distribute the servers.
- Support for load balancing in sidtributed by using connection migration.

**An Erlang High Performance TCP/IP Stack**

**Javier París, Víctor M. Gulías, Alberto Valderruten**

**Introduction**

**Design**

**Optimization Issues**

**Tests**

**Distribution**

# An Erlang High Performance TCP/IP Stack

Javier París    Víctor M. Gulías    Alberto Valderruten

<{paris,gulias,valderruten}@dc.fi.udc.es>