# SCALING ERLANG WEB APPLICATIONS
## 100 TO 100K USERS AT ONE WEB SERVER
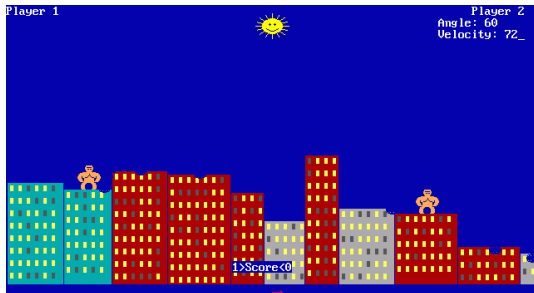
Fernando Benavides (*@elbrujohalcon*)

Inaka Labs

March 20, 2012

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in many dynamic web sites
- Most of them with high scale requirements
- I'll share my experience with you

## HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in many dynamic web sites
- Most of them with high scale requirements
- I'll share my experience with you

## HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in many dynamic web sites
- Most of them with high scale requirements
- I'll share my experience with you

## HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in many dynamic web sites
- Most of them with high scale requirements
- I'll share my experience with you

## HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in many dynamic web sites
- Most of them with high scale requirements
- I'll share my experience with you

# OUTLINE

# OUTLINE

# OUTLINE

# OUTLINE

# OUTLINE

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.

Examples:

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an
*HTTP API* and a component that keeps clients *connected* to the
server for *long periods* of time.

Examples:

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.

Examples:

- Social sites

- IM sites

- Some sites

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.

Examples:

- Social sites
- Chat sites
- Gaming sites

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.
Examples:

- Social sites
- Chat sites
- Sports sites

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.
Examples:

- Social sites
- Chat sites
- Sports sites

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**Description**
Scope

## INTRODUCTION

We will work on the scalability of a *web* project that has an *HTTP API* and a component that keeps clients *connected* to the server for *long periods* of time.
Examples:

- Social sites
- Chat sites
- Sports sites

**Introduction**
The Project
Scaling
Tips and Tricks
Final Words

Description
**Scope**

## SCOPE

*We will try to improve the way we use*

- OTP behaviours
- TCP and HTTP connections
- Underlaying system configurations

*We will **not** deal with*

- Multiple machines/nodes
- Database choices and/or implementations

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**Description**
**Scope**

## SCOPE

*We will try to improve the way we use*

- OTP behaviours
- TCP and HTTP connections
- Underlaying system configurations

*We will **not** deal with*

- Multiple machines/nodes
- Database choices and/or implementations

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

# MATCH STREAM

TODO: General system design graph

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

## REQUIREMENTS

This kind of system requires. . .

- Tons of users
- In two-hour-long bursts
- Real-time updates

Erlang seems to be **the right fit for this**

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

## REQUIREMENTS

This kind of system requires. . .

- Tons of users
- In two-hour-long bursts
- Real-time updates

Erlang seems to be **the right fit for this**

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

## REQUIREMENTS

This kind of system requires. . .

- Tons of users
- In two-hour-long bursts
- Real-time updates

Erlang seems to be **the right fit for this**

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

## REQUIREMENTS

This kind of system requires. . .

- Tons of users
- In two-hour-long bursts
- Real-time updates

Erlang seems to be **the right fit for this**

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

**Idea**
Design
Components

## Requirements

This kind of system requires. . .

- Tons of users
- In two-hour-long bursts
- Real-time updates

Erlang seems to be **the right fit for this**

Introduction
The Project
Scaling
Tips and Tricks
Final Words

Idea
Design
Components

# MATCH STREAM
## GENERAL DESIGN

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

Idea
**Design**
Components

# MATCH STREAM
## ARCHITECTURE

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

Idea
Design
**Components**

# FOR THE CLIENTS

CLIENT_LISTENER

gen_server. Handles a listening TCP port to
receive clients

CLIENT_SUP

supervisor. Supervises client processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb. Listener for HTTP API calls

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

Idea
Design
**Components**

## FOR THE CLIENTS

CLIENT_LISTENER

gen_server. Handles a listening TCP port to receive clients

CLIENT_SUP

supervisor. Supervises client processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb. Listener for HTTP API calls

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

Idea
Design
**Components**

## FOR THE CLIENTS

CLIENT_LISTENER

gen_server. Handles a listening TCP port to receive clients

CLIENT_SUP

supervisor. Supervises client processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb. Listener for HTTP API calls

Introduction
**The Project**
Scaling
Tips and Tricks
Final Words

Idea
Design
**Components**

## FOR THE CLIENTS

CLIENT_LISTENER

gen_server. Handles a listening TCP port to receive clients

CLIENT_SUP

supervisor. Supervises client processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb. Listener for HTTP API calls

## LESSON LEARNED

*Using Erlang to build your system is **not enough** to ensure **scalability***

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
**Blackbox Tests**
**Erlang Tuning**
**Adding Nodes**

# GENERAL RULES
THINGS YOU NEED TO KNOW WHEN SCALING YOUR SYSTEMS

- Start with a system that works

- Automate your clients

- Keep a human watching

- Be patient

# GENERAL RULES
THINGS YOU NEED TO KNOW WHEN SCALING YOUR SYSTEMS

- Start with a system that works

- Automate your clients

- Keep a human watching

- Be patient

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
**Blackbox Tests**
**Erlang Tuning**
**Adding Nodes**

# GENERAL RULES
## THINGS YOU NEED TO KNOW WHEN SCALING YOUR SYSTEMS

- Start with a system that works
- Automate your clients
- Keep a human watching
- Be patient

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
**Blackbox Tests**
**Erlang Tuning**
**Adding Nodes**

# GENERAL RULES
## THINGS YOU NEED TO KNOW WHEN SCALING YOUR SYSTEMS

- Start with a system that works
- Automate your clients
- Keep a human watching
- Be patient

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
Blackbox Tests
Erlang Tuning
Adding Nodes

# STAGE 1
## FINDING THE INITIAL BOUNDARIES

### GOALS

- Test the system as it is
- Find *N* and *C*

### STEPS

- Create the automated testers

- Choose *N* and *C*

- Test the API and long-lived connections independently

- Test both together

- Repeat until you find the highest *N* and *C*

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
Blackbox Tests
Erlang Tuning
Adding Nodes

STAGE 1
FINDING THE INITIAL BOUNDARIES

GOALS

- Test the system as it is
- Find *N* and *C*

STEPS

- Create the automated testers
- Choose *N* and *C*
- Test the API and long-lived connections independently
- Test both together
- Repeat until you find the highest *N* and *C*

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

**Finding The Initial Boundaries**
Blackbox Tests
Erlang Tuning
Adding Nodes

# STAGE 1
## FINDING THE INITIAL BOUNDARIES

## RESULTS



**N = 1024 / C = 4**

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
**Blackbox Tests**
Erlang Tuning
Adding Nodes

# STAGE 2
## BLACKBOX TESTS

### GOALS

- Improve the system environment
- Find the highest *N* and *C* without altering the code

### STEPS

- Check kernel variables

- Check system limits

- Check Erlang VM parameters

- Repeat from Stage 1

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
**Blackbox Tests**
Erlang Tuning
Adding Nodes

# STAGE 2
## BLACKBOX TESTS

GOALS

- Improve the system environment
- Find the highest *N* and *C* without altering the code

STEPS

- Check kernel variables
- Check system limits
- Check Erlang VM parameters
- Repeat from Stage 1

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
**Blackbox Tests**
Erlang Tuning
Adding Nodes

# STAGE 2
## BLACKBOX TESTS

## RESULTS



N = 4096 / C = 4

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
Blackbox Tests
**Erlang Tuning**
Adding Nodes

# STAGE 3
## ERLANG TUNING

GOALS

- Tune up your system
- Find the highest *N* and *C* for one node

STEPS

- Find a problem
- Fix it using the list of *Tips and Tricks*
- If not there, add it
- Repeat from Stage 1

# STAGE 3
## ERLANG TUNING

GOALS

- Tune up your system
- Find the highest *N* and *C* for one node

STEPS

- Find a problem
- Fix it using the list of *Tips and Tricks*
- If not there, add it
- Repeat from Stage 1

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
Blackbox Tests
**Erlang Tuning**
Adding Nodes

# STAGE 3
## ERLANG TUNING

## RESULTS



**N = 65536 / C = 8192**

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
Blackbox Tests
Erlang Tuning
**Adding Nodes**

# STAGE 4
## ADDING NODES

GOALS

- Find the best system topology
- Find *N* and *C* per node

STEPS

- Add a node
  - connected; or
  - independent

- Repeat from Stage 1

# STAGE 4
## ADDING NODES

GOALS

- Find the best system topology
- Find *N* and *C* per node

STEPS

- Add a node
  - connected; or
  - independent
- Repeat from Stage 1

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
Blackbox Tests
Erlang Tuning
**Adding Nodes**

# STAGE 4
## ADDING NODES

## RESULTS



**N ≈ 25000 * 4 ≈ 100000**
**C ≈ 8192 * 4 ≈ 32768**

Introduction
The Project
**Scaling**
Tips and Tricks
Final Words

Finding The Initial Boundaries
Blackbox Tests
Erlang Tuning
**Adding Nodes**

# MATCH STREAM
## FINAL ARCHITECTURE

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

# TIPS AND TRICKS

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
Other Stuff

# OS TWEAKS

## *Kernel Variables*

```
sysctl -w net.ipv4.ip_local_port_range="1024 65535"
sysctl -w net.core.rmem_max=16777216
sysctl -w net.core.wmem_max=16777216
sysctl -w net.ipv4.tcp_rmem="4096 87380 16777216"
sysctl -w net.ipv4.tcp_wmem="4096 65536 16777216"
sysctl -w net.ipv4.tcp_syncookies=1
sysctl -w net.ipv4.tcp_mem="50576    64768    98152"
sysctl -w net.core.netdev_max_backlog=2500
sysctl -w net.netfilter.nf_conntrack_max=1233000
```

## *Open Files Limit*

```
ulimit -n 999999
```

## *Erlang VM tweaks*

+P   Number of Processes

+K   Kernell Polling

-smp   SMP Support

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
Other Stuff

# OS TWEAKS

## *Kernel Variables*

```
sysctl -w net.ipv4.ip_local_port_range="1024 65535"
sysctl -w net.core.rmem_max=16777216
sysctl -w net.core.wmem_max=16777216
sysctl -w net.ipv4.tcp_rmem="4096 87380 16777216"
sysctl -w net.ipv4.tcp_wmem="4096 65536 16777216"
sysctl -w net.ipv4.tcp_syncookies=1
sysctl -w net.ipv4.tcp_mem="50576   64768    98152"
sysctl -w net.core.netdev_max_backlog=2500
sysctl -w net.netfilter.nf_conntrack_max=1233000
```

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
Other Stuff

# OS TWEAKS

## *Kernel Variables*

```
sysctl -w net.ipv4.ip_local_port_range="1024 65535"
sysctl -w net.core.rmem_max=16777216
sysctl -w net.core.wmem_max=16777216
sysctl -w net.ipv4.tcp_rmem="4096 87380 16777216"
sysctl -w net.ipv4.tcp_wmem="4096 65536 16777216"
sysctl -w net.ipv4.tcp_syncookies=1
sysctl -w net.ipv4.tcp_mem="50576    64768    98152"
sysctl -w net.core.netdev_max_backlog=2500
sysctl -w net.netfilter.nf_conntrack_max=1233000
```

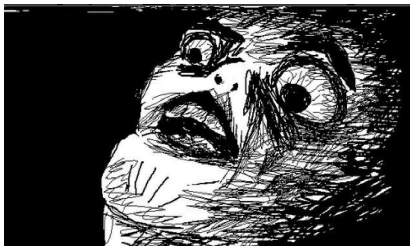## *Open Files Limit*

```
ulimit -n 999999
```

## *Erlang VM tweaks*

$+P$   Number of Processes

$+K$   Kernell Polling

-SMP   SMP Support

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
Other Stuff

# OS TWEAKS

## *Kernel Variables*

```
sysctl -w net.ipv4.ip_local_port_range="1024 65535"
sysctl -w net.core.rmem_max=16777216
sysctl -w net.core.wmem_max=16777216
sysctl -w net.ipv4.tcp_rmem="4096 87380 16777216"
sysctl -w net.ipv4.tcp_wmem="4096 65536 16777216"
sysctl -w net.ipv4.tcp_syncookies=1
sysctl -w net.ipv4.tcp_mem="50576    64768    98152"
sysctl -w net.core.netdev_max_backlog=2500
sysctl -w net.netfilter.nf_conntrack_max=1233000
```

## *Open Files Limit*

```
ulimit -n 999999
```

## *Erlang VM tweaks*

$+P$  Number of Processes

$+K$  Kernell Polling

-SMP  SMP Support

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## CONNECTION TWEAKS

### BACKLOG

- Allow more concurrent connections
- Remember HTTP *runs on* TCP

### CONNECTIONS

- Don't use just one of them
- Check inbound and outbound connections

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## CONNECTION TWEAKS

### BACKLOG

- Allow more concurrent connections
- Remember HTTP *runs on* TCP

### CONNECTIONS

- Don't use just one of them
- Check inbound and outbound connections

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## GEN_EVENT

### SUP_HANDLER

- Don't use it
- Monitor the processes instead

### LONG DELIVERY QUEUES

- Use *repeaters*

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## GEN_EVENT

SUP_HANDLER

- Don't use it
- Monitor the processes instead

LONG DELIVERY QUEUES

- Use *repeaters*

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## GEN_SERVER

### CALL TIMEOUTS
Remember `gen_server:reply/2`

### MEMORY FOOTPRINT
Remember `hibernate`

### LONG INIT/1
Use 0 timeout

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
**OTP**
Other Stuff

## GEN_SERVER

### CALL TIMEOUTS
Remember gen_server:reply/2

### MEMORY FOOTPRINT
Remember hibernate

### LONG INIT/1
Use 0 timeout

**Introduction**
**The Project**
**Scaling**
**Tips and Tricks**
**Final Words**

**TCP Tunning**
**OTP**
**Other Stuff**

## GEN_SERVER

CALL TIMEOUTS
Remember gen_server:reply/2

MEMORY FOOTPRINT
Remember hibernate

LONG INIT/1
Use 0 timeout

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
**OTP**
Other Stuff

## SUPERVISORS

- Sometimes `simple_one_for_one` supervisors get overburdened because they have too many children
- Try a supervisor hierarchy with several managers below the main supervisor
- Turn `supervisor:start_child/2` calls into something like

```
supervisor:start_child(
    list_to_atom("module-name_" ++
                            integer_to_list(random:uniform(#ofSupervisors))).
```

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
**Other Stuff**

# OTHER PROCESSES

### TIMERS

- Don't use the timer module
- Use erlang:send_after

### LOGGING

- Don't log too much
- Use a good logging system

### REGISTRATION

- Sometimes it's better to register processes instead of keeping track of their pids manually
- You can always register processes both locally and globally

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
OTP
**Other Stuff**

# OTHER PROCESSES

### TIMERS

- Don't use the `timer` module
- Use `erlang:send_after`

### LOGGING

- Don't log too much
- Use a good logging system

### REGISTRATION

- Sometimes it's better to register processes instead of keeping track of their pids manually
- You can always register processes both locally and globally

Introduction
The Project
Scaling
**Tips and Tricks**
Final Words

**TCP Tunning**
**OTP**
**Other Stuff**

# OTHER PROCESSES

### TIMERS

- Don't use the `timer` module
- Use `erlang:send_after`

### LOGGING

- Don't log too much
- Use a good logging system

### REGISTRATION

- Sometimes it's better to register processes instead of keeping track of their pids manually
- You can always register processes <span style="color:red">both</span> locally and globally

Introduction
The Project
Scaling
Tips and Tricks
Final Words

Summary
What's next?
Questions

# SUMMARY

- This is an iterative process
- It worked awesomely for us in both experimental and real-life systems
- It's no silver bullet
- The list of *Tips and Tricks* grows constantly over time

Introduction
The Project
Scaling
Tips and Tricks
Final Words

Summary
What's next?
Questions

## Summary

- This is an iterative process
- It worked awesomely for us in both experimental and real-life systems
- It's no silver bullet
- The list of *Tips and Tricks* grows constantly over time

Introduction
The Project
Scaling
Tips and Tricks
Final Words

Summary
What's next?
Questions

# SUMMARY

- This is an iterative process
- It worked awesomely for us in both experimental and real-life systems
- It's no silver bullet
- The list of *Tips and Tricks* grows constantly over time

Introduction
The Project
Scaling
Tips and Tricks
Final Words

**Summary**
What's next?
Questions

# SUMMARY

- This is an iterative process
- It worked awesomely for us in both experimental and real-life systems
- It's no silver bullet
- The list of *Tips and Tricks* grows constantly over time

Introduction
The Project
Scaling
Tips and Tricks
**Final Words**

Summary
**What's next?**
Questions

# SCALING TOPICS
## THAT WEREN'T COVERED ON THIS PRESENTATION

- Adding nodes
- Choosing databases
- System specific improvements
- Measuring tools

Introduction
The Project
Scaling
Tips and Tricks
Final Words

Summary
What's next?
Questions

# QUESTIONS

# Thanks!