# SCALING ERLANG WEB APPLICATIONS
## 100 TO 100K USERS AT ONE WEB SERVER

Fernando Benavides (*@elbrujohalcon*)

Inaka Labs

March 21, 2012

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers

- I'll show you how I make them scale

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers
  - Many of them with real-time updates
  - And in the right kinds of environments
- I'll show you how I make them scale



**Fernando Benavides (*@elbrujohalcon*)**     **Scaling Erlang Web Applications**

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers
  - Many of them with real-time updates
  - Most of them with high scale requirements
- I'll show you how I make them scale

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers
  - Many of them with real-time updates
  - Most of them with high scale requirements
- I'll show you how I make them scale

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers
  - Many of them with real-time updates
  - Most of them with high scale requirements
- I'll show you how I make them scale

# HELLO WORLD!

- I'm a developer since I was 10
- I'm an Erlang developer since 2008
- I've worked in several dynamic web servers
  - Many of them with real-time updates
  - Most of them with high scale requirements
- I'll show you how I make them scale

**Introduction**
**Match Stream**
**Scaling**
**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**
**Match Stream**
**Scaling**
**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**
**Match Stream**
**Scaling**
**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 apps
- MMRPG sites

**Introduction**
**Match Stream**
**Scaling**
**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

- HTTP sites

**Introduction**
Match Stream
Scaling
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**
**Match Stream**
**Scaling**
**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**
Match Stream
Scaling
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**
Match Stream
Scaling
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**
Match Stream
Scaling
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
**Match Stream**
Scaling
Final Words

**Idea**
Design
Components

# MATCH STREAM

TODO: General system design graph

Introduction
**Match Stream**
Scaling
Final Words

**Idea**
Design
Components

# MATCH STREAM
### REQUIREMENTS

### SYSTEM CHALLENGES

- Tons of concurrent users
- Two-hour-long bursts of connections followed by long periods of inactivity
- Real-time updates

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

Introduction
**Match Stream**
Scaling
Final Words

**Idea**
Design
Components

# MATCH STREAM
## REQUIREMENTS

SYSTEM CHALLENGES

- Tons of concurrent users
- Two-hour-long bursts of connections followed by long periods of inactivity
- Real-time updates

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

Introduction
**Match Stream**
Scaling
Final Words

**Idea**
Design
Components

# MATCH STREAM
## REQUIREMENTS

SYSTEM CHALLENGES

- Tons of concurrent users
- Two-hour-long bursts of connections followed by long periods of inactivity
- Real-time updates

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

Introduction
**Match Stream**
Scaling
Final Words

**Idea**
Design
Components

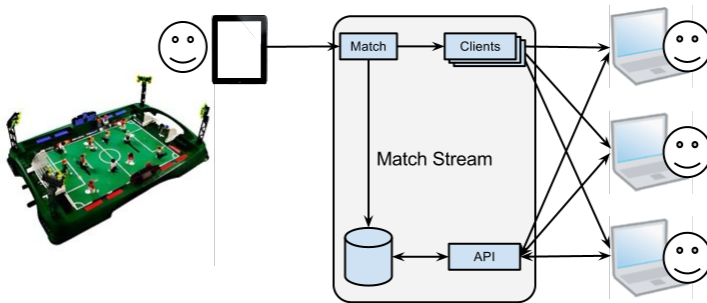# MATCH STREAM
### REQUIREMENTS

SYSTEM CHALLENGES

- Tons of concurrent users
- Two-hour-long bursts of connections followed by long periods of inactivity
- Real-time updates

Erlang seems to be **the right fit for this**
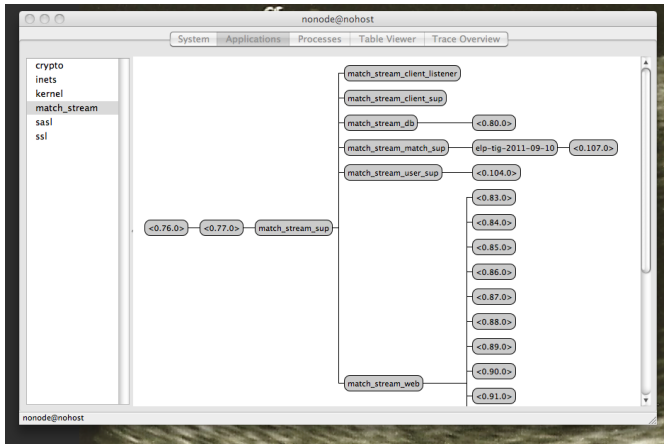
Introduction
**Match Stream**
Scaling
Final Words

Idea
**Design**
Components

# MATCH STREAM
## GENERAL DESIGN

Introduction
**Match Stream**
Scaling
Final Words

**Idea**
**Design**
Components

# MATCH STREAM
## ARCHITECTURE

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

MATCH STREAM
CLIENT-HANDLING COMPONENTS

CLIENT_LISTENER

gen_server. Listens on a TCP port to receive client connections

CLIENT_SUP

supervisor. Supervises connection processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb server. Listens for HTTP API calls

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

MATCH STREAM
CLIENT-HANDLING COMPONENTS

CLIENT_LISTENER

gen_server. Listens on a TCP port to receive client connections

CLIENT_SUP

supervisor. Supervises connection processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb server. Listens for HTTP API calls

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

# MATCH STREAM
## CLIENT-HANDLING COMPONENTS

CLIENT_LISTENER

> `gen_server`. Listens on a TCP port to receive client connections

CLIENT_SUP

> `supervisor`. Supervises connection processes

USER_SUP

> `supervisor`. Supervises user processes

WEB

> `mochiweb server`. Listens for HTTP API calls

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

# MATCH STREAM
## CLIENT-HANDLING COMPONENTS

CLIENT_LISTENER

gen_server. Listens on a TCP port to receive client connections

CLIENT_SUP

supervisor. Supervises connection processes

USER_SUP

supervisor. Supervises user processes

WEB

mochiweb server. Listens for HTTP API calls

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

MATCH STREAM
DB AND WATCHER COMPONENTS

DB

`gen_server`. Handles a connection to the DB

MATCH_SUP

`supervisor`. Supervises match processes

Introduction
**Match Stream**
Scaling
Final Words

Idea
Design
**Components**

# MATCH STREAM
## DB AND WATCHER COMPONENTS

DB

gen_server. Handles a connection to the DB

MATCH_SUP

supervisor. Supervises match processes

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# LESSON LEARNED

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

# STAGE 1
## TESTING THE SYSTEM AS IT IS

### GOALS

- Find how much the system can handle

### STEPS

- Create automated testers

- Start the system on a *clean* machine

- Test repeatedly adjusting the number of connections

- Have a human trying the system himself

Introduction
Match Stream
**Scaling**
Final Words

**Stage 1: The Original System**
Stage 2: OS Tuning
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 1
## TESTING THE SYSTEM AS IT IS

GOALS
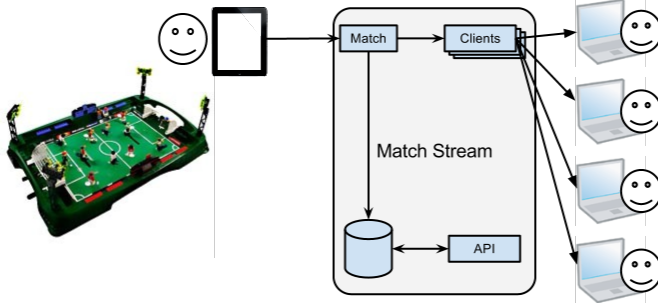
- Find how much the system can handle

STEPS

- Create automated testers
- Start the system on a *clean* machine
- Test repeatedly adjusting the number of connections
- Have a human trying the system himself

Introduction
Match Stream
**Scaling**
Final Words

**Stage 1: The Original System**
Stage 2: OS Tuning
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 1
## TESTING THE SYSTEM AS IT IS

## RESULTS



N = 1024 / C = 4

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
**Stage 2: OS Tuning**
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 2
## IMPROVING THE ENVIRONMENT

### GOALS

- Improve the system environment without altering the code

### SETTINGS TO TUNE UP

- Concurrent TCP connections
- Open files limit
- TCP memory size
- TCP memory allocation
- Erlang vm startup parameters

# STAGE 2
## IMPROVING THE ENVIRONMENT

### GOALS

- Improve the system environment without altering the code

### SETTINGS TO TUNE UP

- Concurrent TCP connections
- Open files limit
- TCP backlog size
- TCP memory allocation
- Erlang VM startup parameters

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
**Stage 2: OS Tuning**
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 2
## IMPROVING THE ENVIRONMENT

GOALS

- Improve the system environment without altering the code

SETTINGS TO TUNE UP

- Concurrent TCP connections

- Open files limit

- TCP backlog size

- TCP memory allocation

- Erlang VM startup parameters

## STAGE 2
### IMPROVING THE ENVIRONMENT

GOALS

- Improve the system environment without altering the code

SETTINGS TO TUNE UP

- Concurrent TCP connections
- Open files limit
- TCP backlog size
- TCP memory allocation
- Erlang VM startup parameters

# STAGE 2
## IMPROVING THE ENVIRONMENT

GOALS

- Improve the system environment without altering the code

SETTINGS TO TUNE UP

- Concurrent TCP connections
- Open files limit
- TCP backlog size
- TCP memory allocation
- Erlang VM startup parameters

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
**Stage 2: OS Tuning**
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 2
## IMPROVING THE ENVIRONMENT

GOALS

- Improve the system environment without altering the code

SETTINGS TO TUNE UP
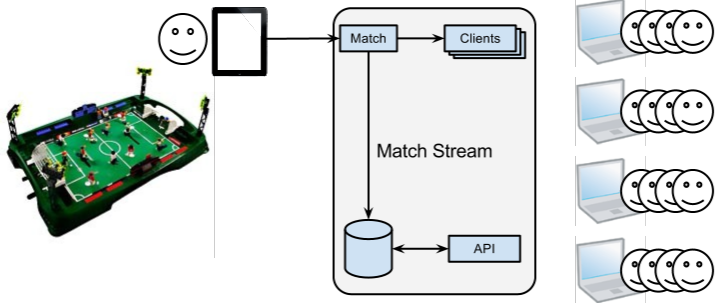
- Concurrent TCP connections
- Open files limit
- TCP backlog size
- TCP memory allocation
- Erlang VM startup parameters

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
**Stage 2: OS Tuning**
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 2
## IMPROVING THE ENVIRONMENT

## RESULTS



**N = 4096 / C = 4**

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3
## IMPROVING MATCH STREAM

### GOALS

- Tune up the system 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
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3
## IMPROVING MATCH STREAM

GOALS

- Tune up the system for <span style="color:red">one node</span>

STEPS

- Find a problem
- Fix it using the list of *Tips and Tricks*
- If not there, add it
- Repeat from <span style="color:red">Stage 1</span>

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.1
## 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
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# Stage 3.1
## Connection Tweaks

### Backlog

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

### Connections

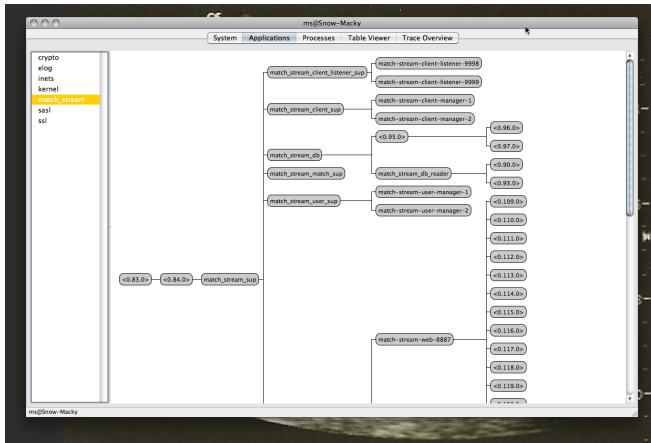- Don't use just one of them
- Check inbound and outbound connections

Introduction
Match Stream
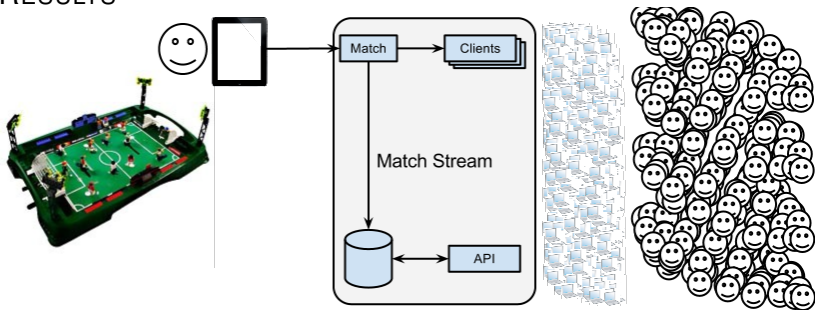**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.1
## CONNECTION TWEAKS

## SYSTEM ARCHITECTURE

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.1
## CONNECTION TWEAKS

## RESULTS



**N = 65536 / C = 8192**

# STAGE 3.2
GEN_EVENT

#### SUP_HANDLER

- Don't use it
- Monitor the processes instead

#### LONG DELIVERY QUEUES

- Use *repeaters*

# STAGE 3.2
GEN_EVENT

SUP_HANDLER

- Don't use it
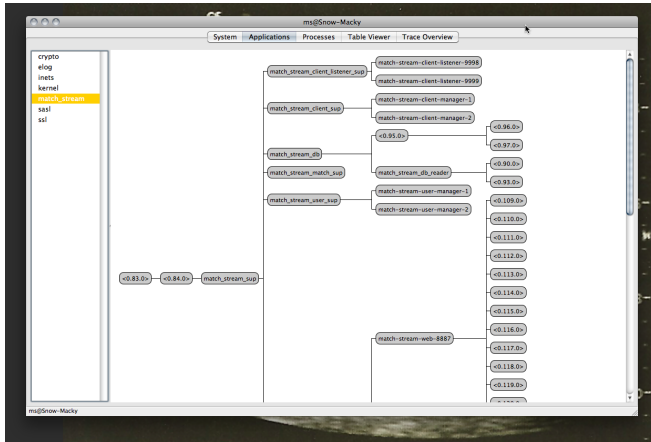- Monitor the processes instead

LONG DELIVERY QUEUES

- Use *repeaters*

Introduction
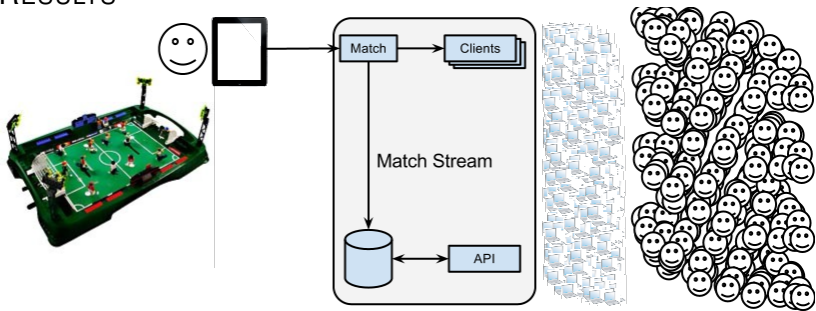Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.2
GEN_EVENT

## SYSTEM ARCHITECTURE

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.2
## GEN_EVENT

## RESULTS



**N = 65536 / C = 8192**

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.3
GEN_SERVER

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

MEMORY FOOTPRINT
Remember hibernate

LONG INIT/1
Use 0 timeout

# STAGE 3.3
GEN_SERVER

CALL TIMEOUTS

Remember gen_server:reply/2
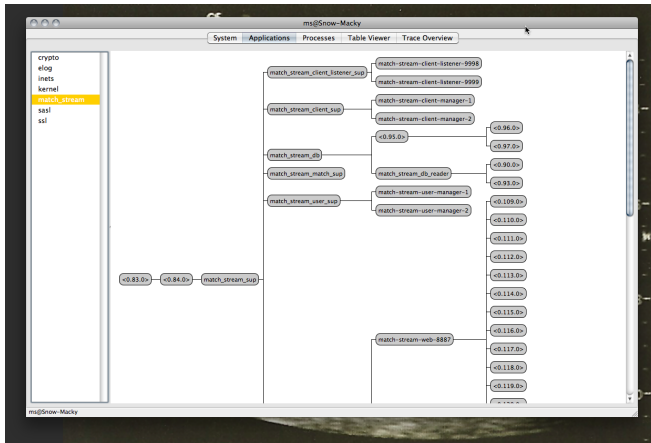
MEMORY FOOTPRINT

Remember hibernate

LONG INIT/1

Use 0 timeout

Introduction    Stage 1: The Original System
Match Stream    Stage 2: OS Tuning
**Scaling**    **Stage 3: Erlang Tuning**
Final Words    Stage 4: Multi-Node Tuning

# STAGE 3.3
GEN_SERVER

CALL TIMEOUTS

Remember `gen_server:reply/2`

MEMORY FOOTPRINT

Remember `hibernate`

LONG INIT/1

Use 0 timeout

Introduction
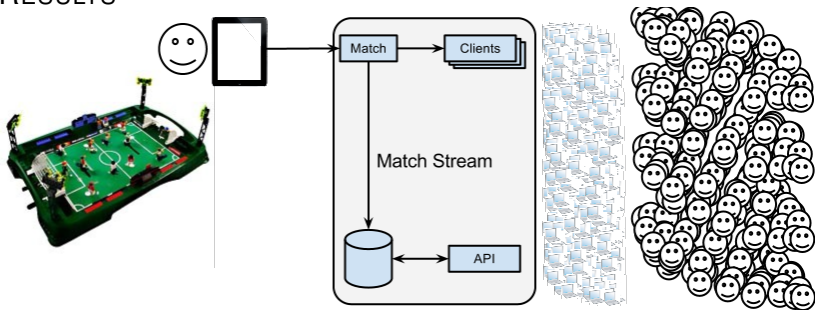Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.3
## GEN_SERVER

## SYSTEM ARCHITECTURE

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.3
GEN_SERVER

## RESULTS



**N = 65536 / C = 8192**

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.4
## 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
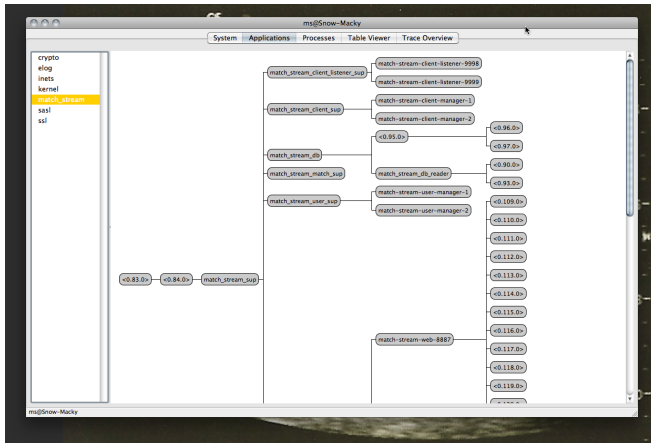- Turn `supervisor:start_child/2` calls into something like

```
supervisor:start_child(
  list_to_atom("module-name_" ++
                    integer_to_list(random:uniform(#ofSupervisors))).
```
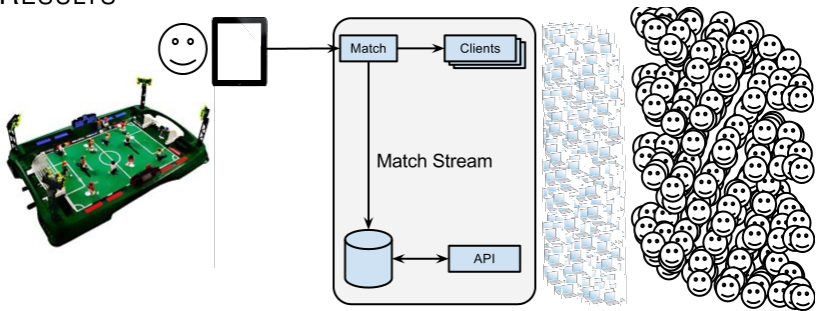
Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.4
## SUPERVISORS

## SYSTEM ARCHITECTURE

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.4
## SUPERVISORS

## RESULTS



**N = 65536 / C = 8192**

# STAGE 3.5
## 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
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.5
## 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
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
**Stage 3: Erlang Tuning**
Stage 4: Multi-Node Tuning

# STAGE 3.5
## OTHER PROCESSES

### TIMERS

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

### LOGGING

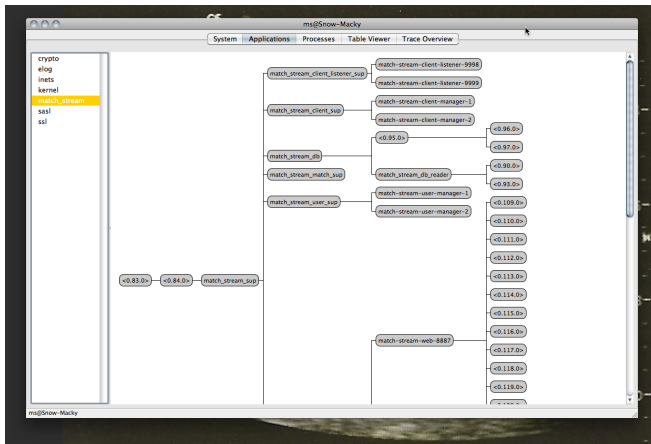- Don't log too much
- Use a good logging system

### REGISTRATION

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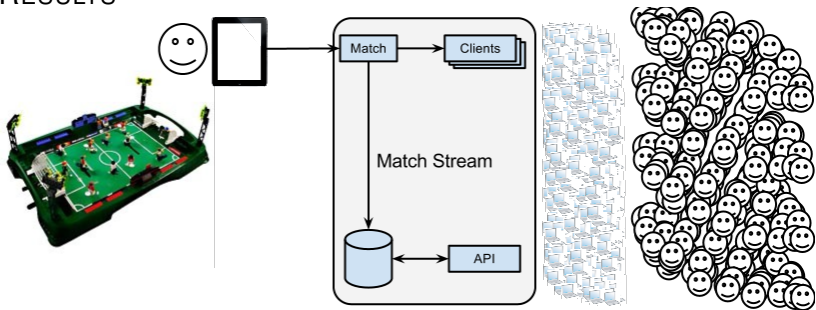- You can always register processes both locally and globally

Introduction
Match Stream
Scaling
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 3.5
## OTHER PROCESSES

## SYSTEM ARCHITECTURE

Introduction
Match Stream
Scaling
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
Stage 4: Multi-Node Tuning

# STAGE 3.5
## OTHER PROCESSES

## RESULTS



**N = 65536 / C = 8192**

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
**Stage 4: Multi-Node Tuning**

# STAGE 4
## ADDING NODES

### GOALS

- Find the best system topology

### STEPS

- Prepare the system to run in more than one node

- Decide if nodes should be connected or independent

- Decide if nodes should be on the same machine or not

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
**Stage 4: Multi-Node Tuning**

# STAGE 4
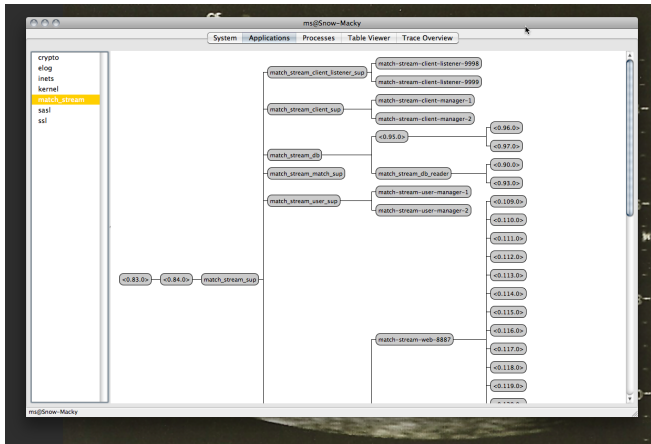## ADDING NODES

GOALS

- Find the best system topology

STEPS

- Prepare the system to run in more than one node
- Decide if nodes should be connected or independent
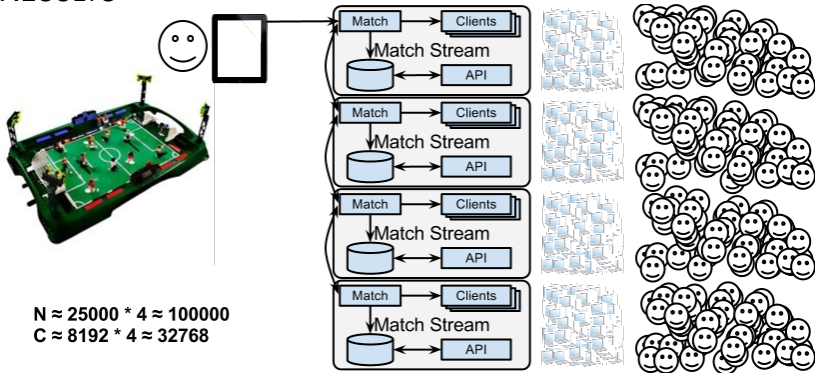- Decide if nodes should be on the same machine or not

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
**Stage 4: Multi-Node Tuning**

# STAGE 4
## ADDING NODES

Introduction
Match Stream
**Scaling**
Final Words

Stage 1: The Original System
Stage 2: OS Tuning
Stage 3: Erlang Tuning
**Stage 4: Multi-Node Tuning**

# STAGE 4
## ADDING NODES

## RESULTS



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

Introduction
Match Stream
Scaling
**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
Match Stream
Scaling
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
Match Stream
Scaling
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
Match Stream
Scaling
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
Match Stream
Scaling
Final Words

Summary
What's next?
Questions

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

- Managing many nodes
- Choosing databases
- System specific improvements
- Measuring tools

Introduction
Match Stream
Scaling
**Final Words**

Summary
What's next?
**Questions**

# QUESTIONS

# Thanks!