

/

Building a Multi-Master, Replicated Redis Clone with Erlang/OTP

Chad DePue

Founder of **Inaka**
CTO of **WhisperText**

Enjoy Erlang
Love Redis
@chaddepue

Who is Inaka?

We are an international team of 25 developers, building high performance iOS/Android apps, using Ruby, Erlang, Elixir, Node.

What is Whisper?





Top 10 Social Network

What is Edis?

Redis - a C-based fast in-memory,
disk-backed key/value database

Edis - an Erlang-based, leveldb-
backed key/value store that speaks the
Redis protocol

I love Redis!

Edis - an Erlang-based server that...

- Uses `gen_tcp`
- Uses `gen_fsm`
- Uses LevelDB
- Implements Redis command set
- Respects Redis algorithms

Why is Redis Great?

- Speed
- Expressivity of the command set
- Ease of Deployment

Command Group

Key/Value

Hashes

Lists

Sets

Sorted Sets

Publish/Subscribe

Transactions

Selected Commands

SET/GET

HSETNX

RPOP/LPOP

SUNION/SPOP

ZADD/ZRANGE

SUBSCRIBE/PUBLISH

MULTI/EXEC

RPOPLPUSH

work_queue

1 , 2 , 7 , 10

RPOPLPUSH

work_queue

1 , 2 , 7, 10

run_queue

18

RPOPLPUSH

work_queue

1 , 2 , 7, 10

run_queue

18



RPOPLPUSH

work_queue

1, 2, 7

run_queue

10, 18



RPOPLPUSH

work_queue

1, 2, 7

run_queue

10, 18



Atomic

Redis Design Decisions

- Data must fit in-memory
- Master-slave model
- Scripting Blocks the Thread

Edis Design Decisions

- Disk-backed by default
- Pluggable DB
- Master-Slave or Master-Master
- Extensible with Erlang

Edis Design Decisions

It's important that Edis respects Redis's goals of algorithmic complexity.

If a Redis command is $O(\log(n))$, Edis will have the same $O()$.*

* Except for ZSETS - We don't yet have skiplists in Erlang.

Edis Use Cases

- Smart Redis Proxy
- Handy Interface on top of existing data
- Reference Protocol Implementation

We're going to ...

- Learn Erlang in 45 seconds
- Review Edis Architecture
- Trace a SET command
- Learn about Replication and MM
- Demo Multi-Master (time permitting)

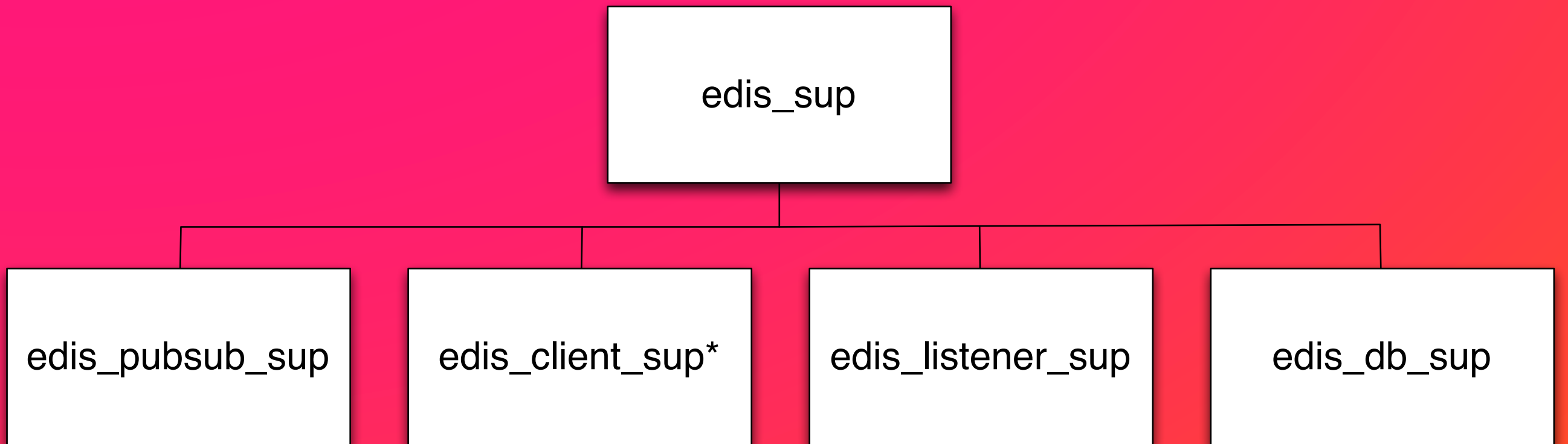
Key terms in Erlang/ OTP in 45 seconds

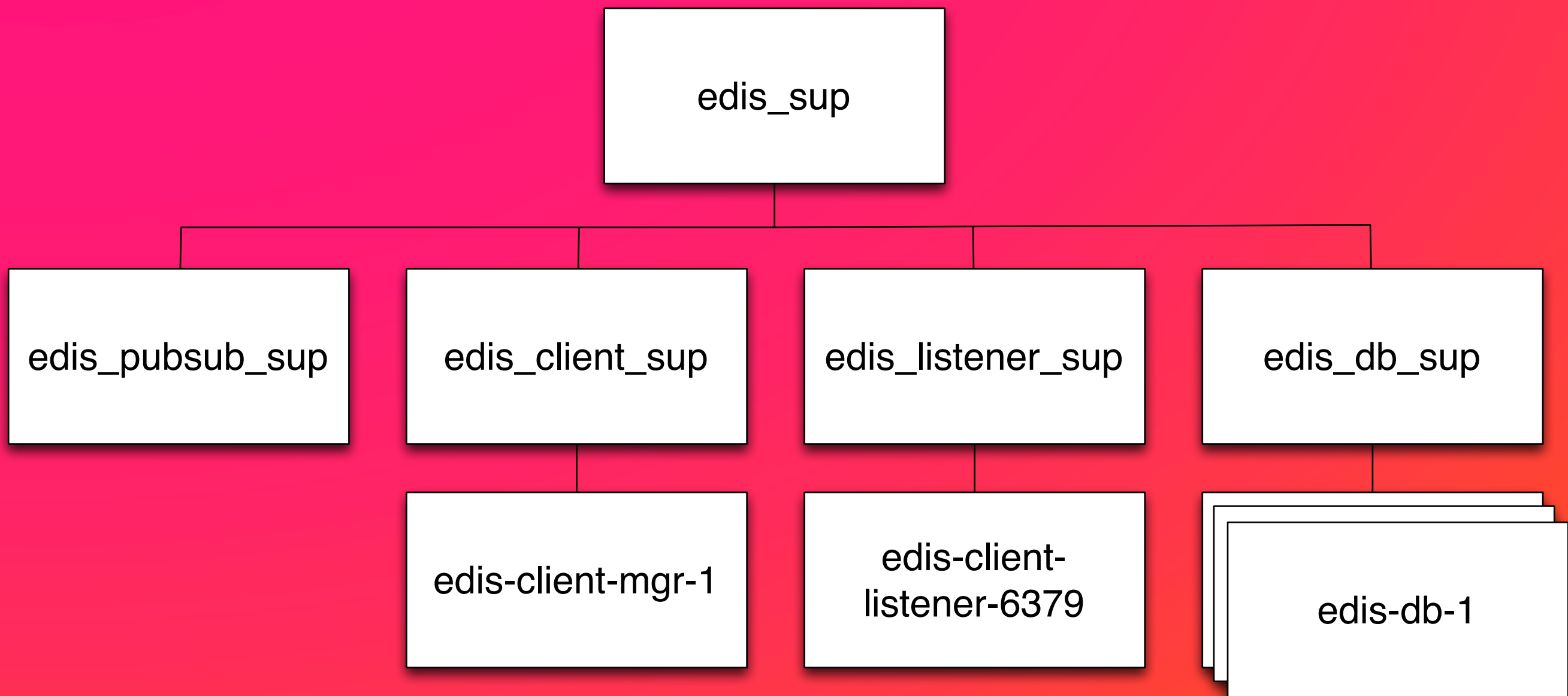
- pattern matching
- variable immutability
- lightweight processes + messages
- semicolon; comma, period.
- OTP provides:
 - gen_server, gen_fsm
 - supervisors (and more)

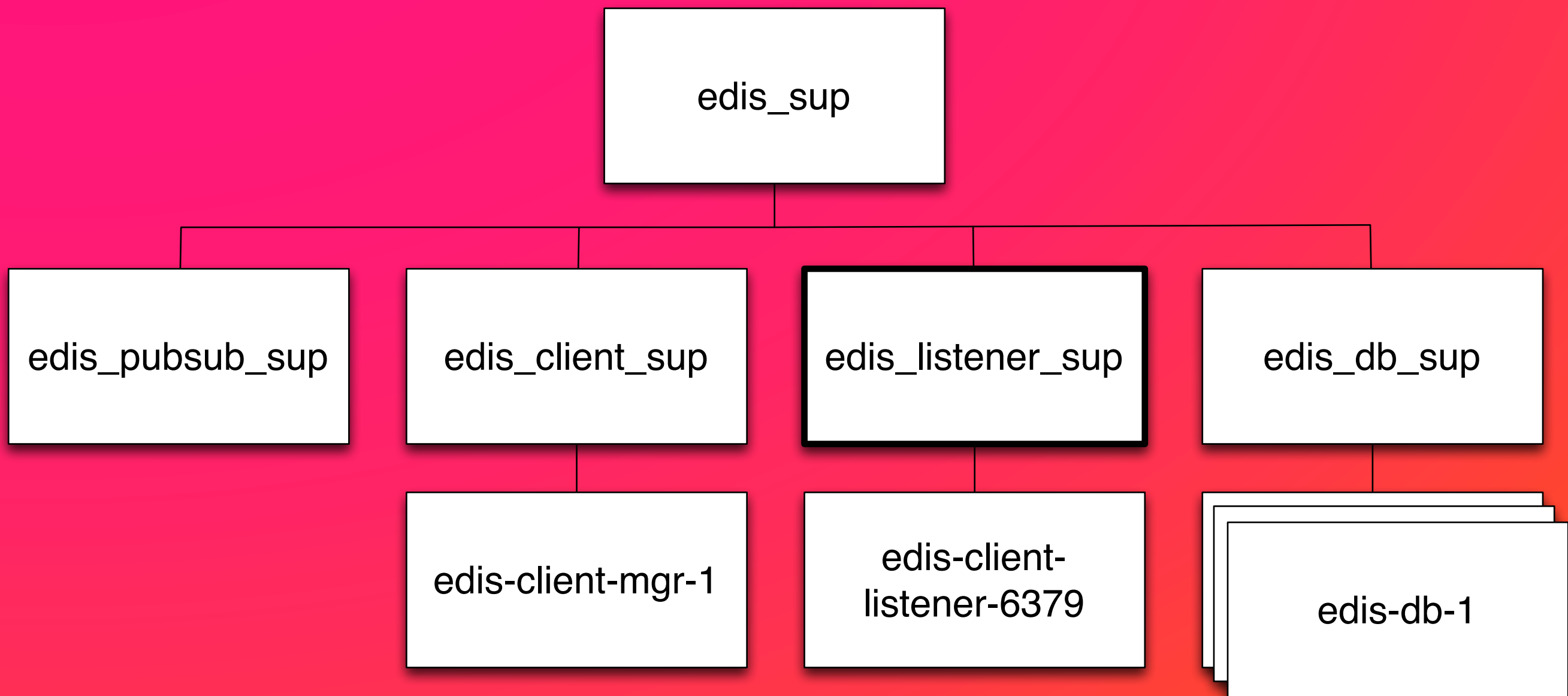
Get All 3 Major Erlang Books

Armstrong Book	(Become a Berliever)
Erlang Programming	(Architecture)
Erlang and OTP in Action	(Practical)

Structure of the App

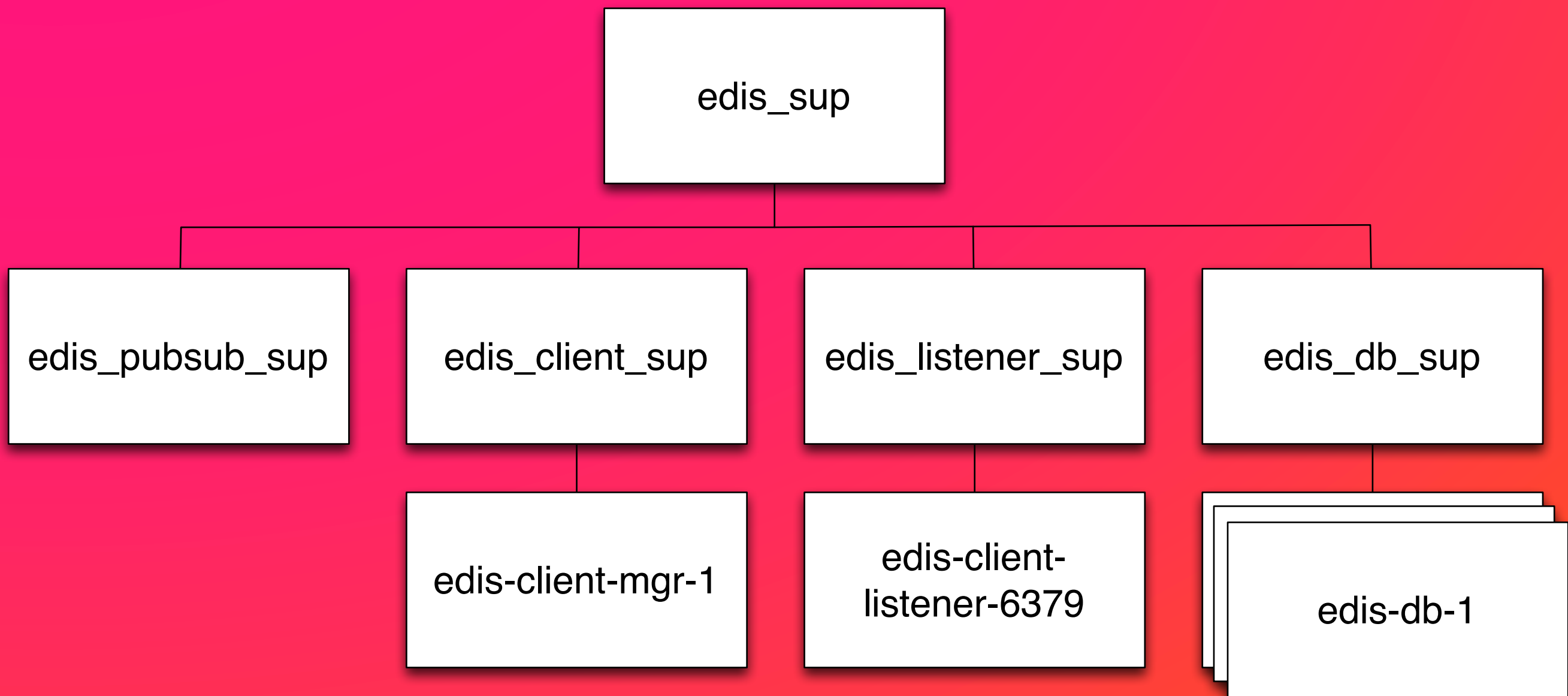


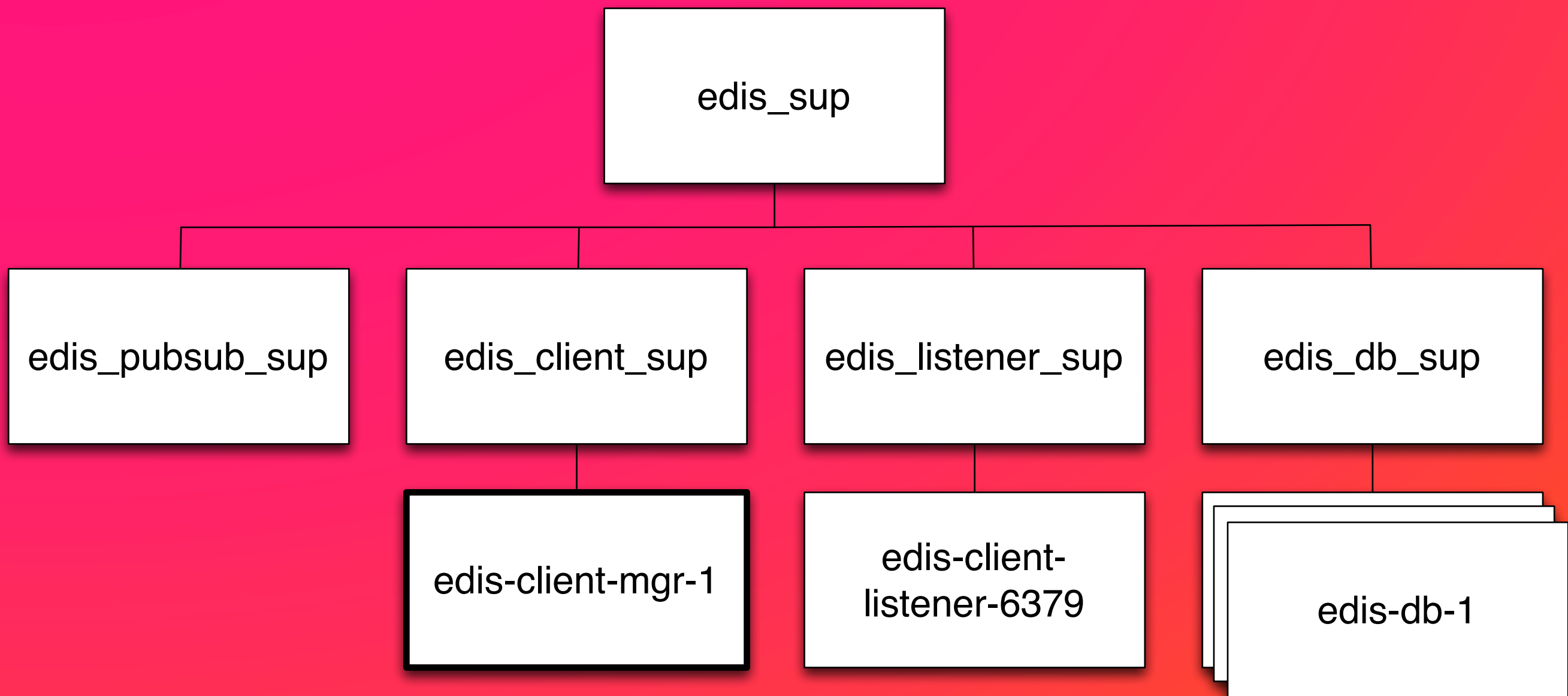


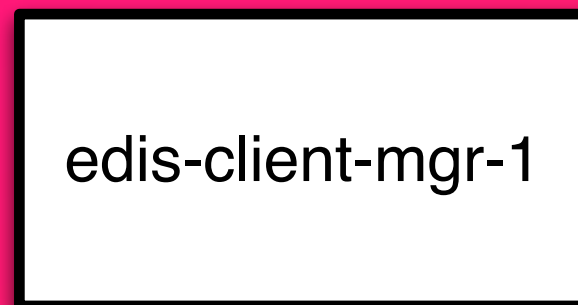


```
37 init([]) ->
38   {MinPort, MaxPort} = edis_config:get(listener_port_range),
39   Listeners =
40     [{list_to_atom("edis-listener-" ++ integer_to_list(I)),
41      {edis_listener, start_link, [I]}, permanent, brutal_kill,
42      worker, [edis_listener]}
43     || I <- lists:seq(MinPort, MaxPort)],
44   {ok, {{one_for_one, 5, 10}, Listeners}}.
```

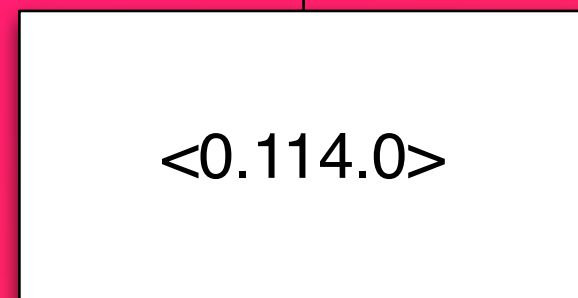
```
37 init([]) ->
38   {MinPort, MaxPort} = edis_config:get(listener_port_range),
39   Listeners =
40     [{list_to_atom("edis-listener-" ++ integer_to_list(I)),
41      {edis_listener, start_link, [I]}, permanent, brutal_kill,
42      worker, [edis_listener]}
43     || I <- lists:seq(MinPort, MaxPort)],
44   {ok, {{one_for_one, 5, 10}, Listeners}}.
```



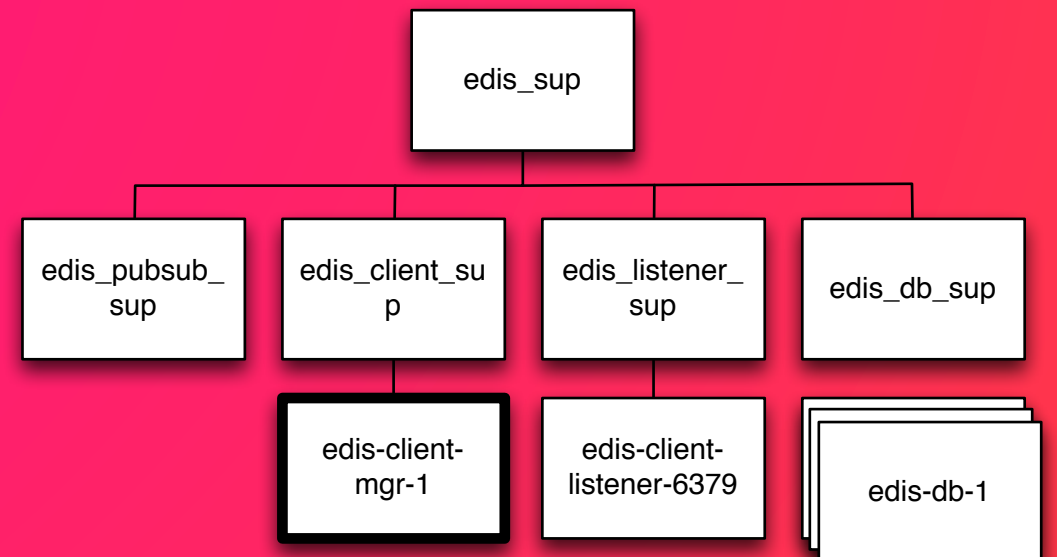
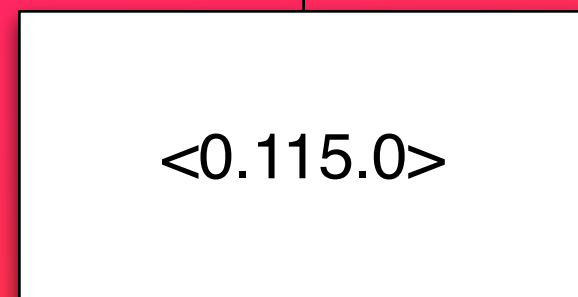




edis_client



edis_command_runner




```
64 socket({socket_ready, Socket}, State) ->
65     % Now we own the socket
66     PeerPort = inet:peername(Socket),
67
68     ok = inet:setopts(Socket, [{active, once},
69                               {packet, line}, binary]),
70     _ = erlang:process_flag(trap_exit, true),
71     {ok, CmdRunner} = edis_command_runner:start_link(Socket),
72     {next_state, command_start,
73       State#state{socket          = Socket,
74                     peerport      = PeerPort,
75                     command_runner = CmdRunner}, hibernate};
```

edis_client.erl

```
64 socket({socket_ready, Socket}, State) ->
65     % Now we own the socket
66     PeerPort = inet:peername(Socket),
67
68     ok = inet:setopts(Socket, [{active, once},
69                                {packet, line}, binary]),
70     _ = erlang:process_flag(trap_exit, true),
71     {ok, CmdRunner} = edis_command_runner:start_link(Socket),
72     {next_state, command_start,
73      State#state{socket          = Socket,
74                    peerport      = PeerPort,
75                    command_runner = CmdRunner}, hibernate};
```

edis_client.erl

```
64 socket({socket_ready, Socket}, State) ->
```

```
{next_state, StateName, State};
```

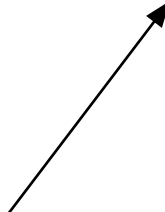
```
72 {next_state, command_start,  
73     State#state{socket      = Socket,  
74                     peerport = PeerPort,  
75                     command_runner = CmdRunner}, hibernate};
```

edis_client.erl

```
64 socket({socket_ready, Socket}, State) ->
```

```
{next_state, StateName, State};
```

```
72 {next_state, command_start,  
73     State#state{socket      = Socket,  
74                     peerport = PeerPort,  
75                     command_runner = CmdRunner}, hibernate};
```



edis_client.erl

```
gen_fsm:send_event(<Pid>,message).
```



```
edis_client:StateName(message,State).
```

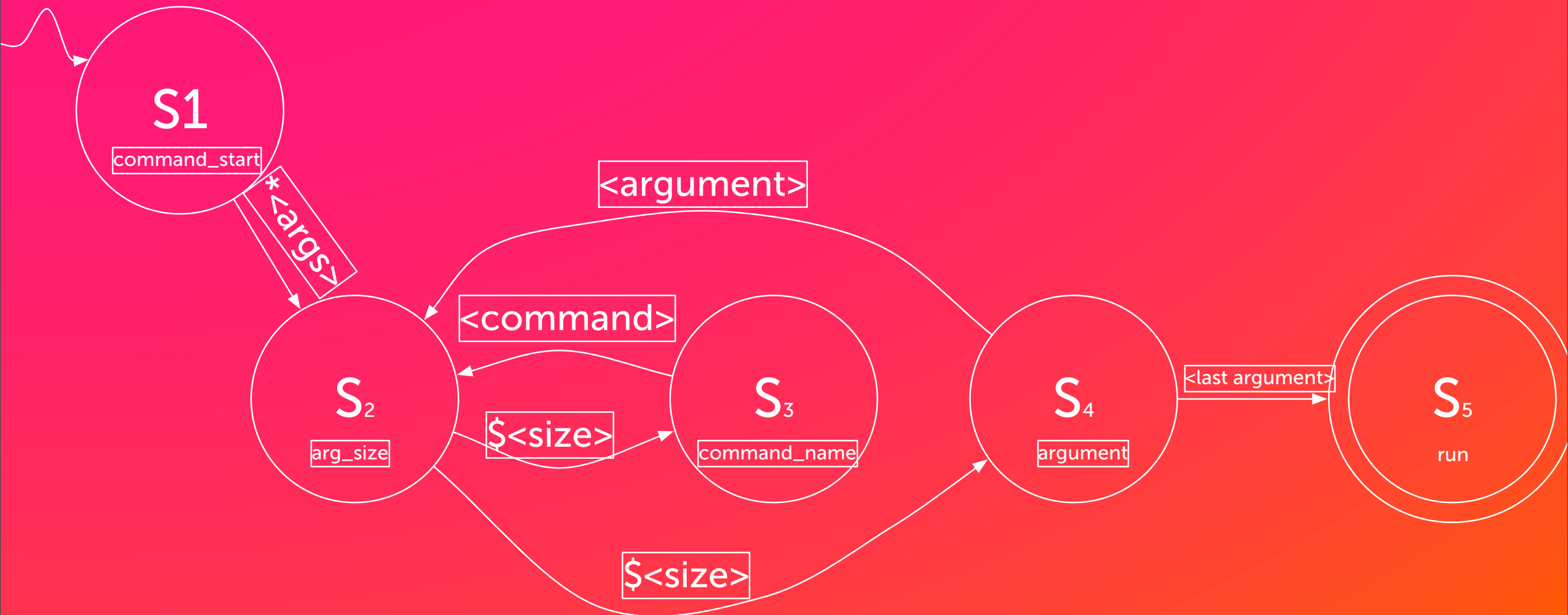
edis_client.erl

```
edis_client:command_start({data,Data},State).
```

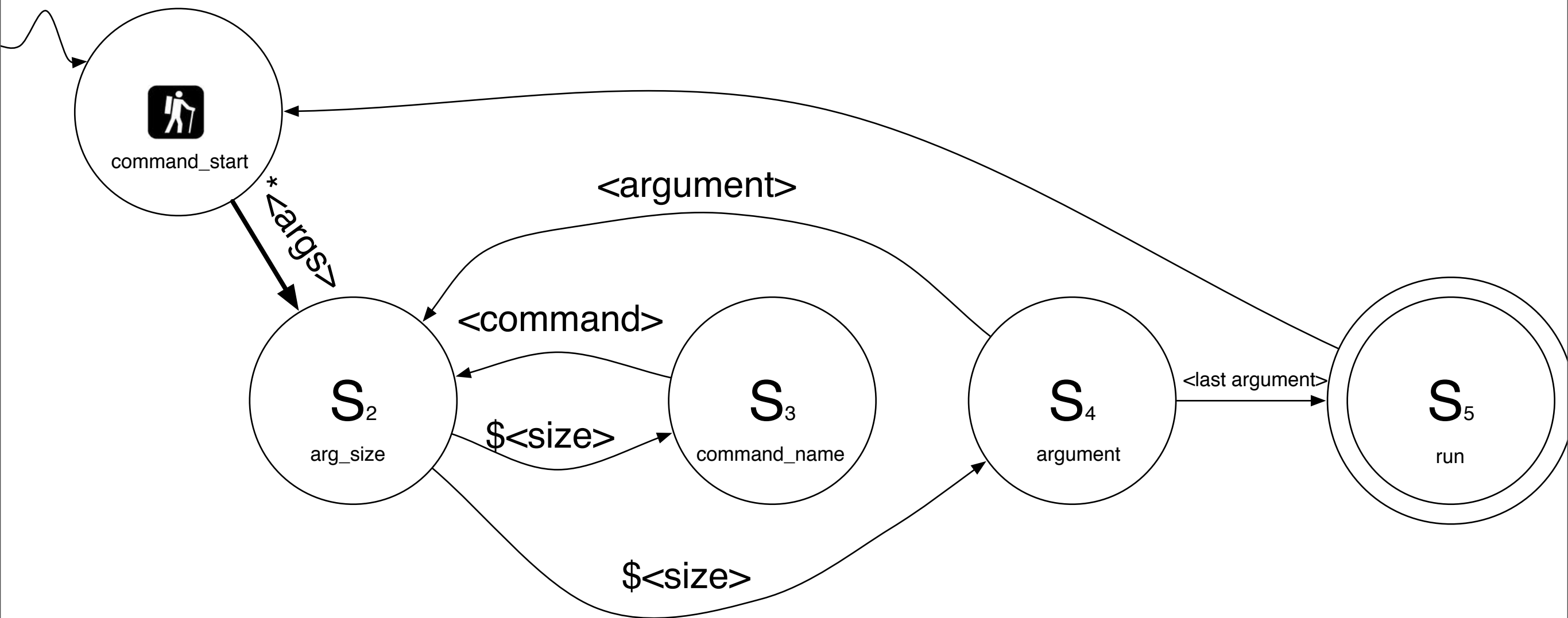
Redis Protocol



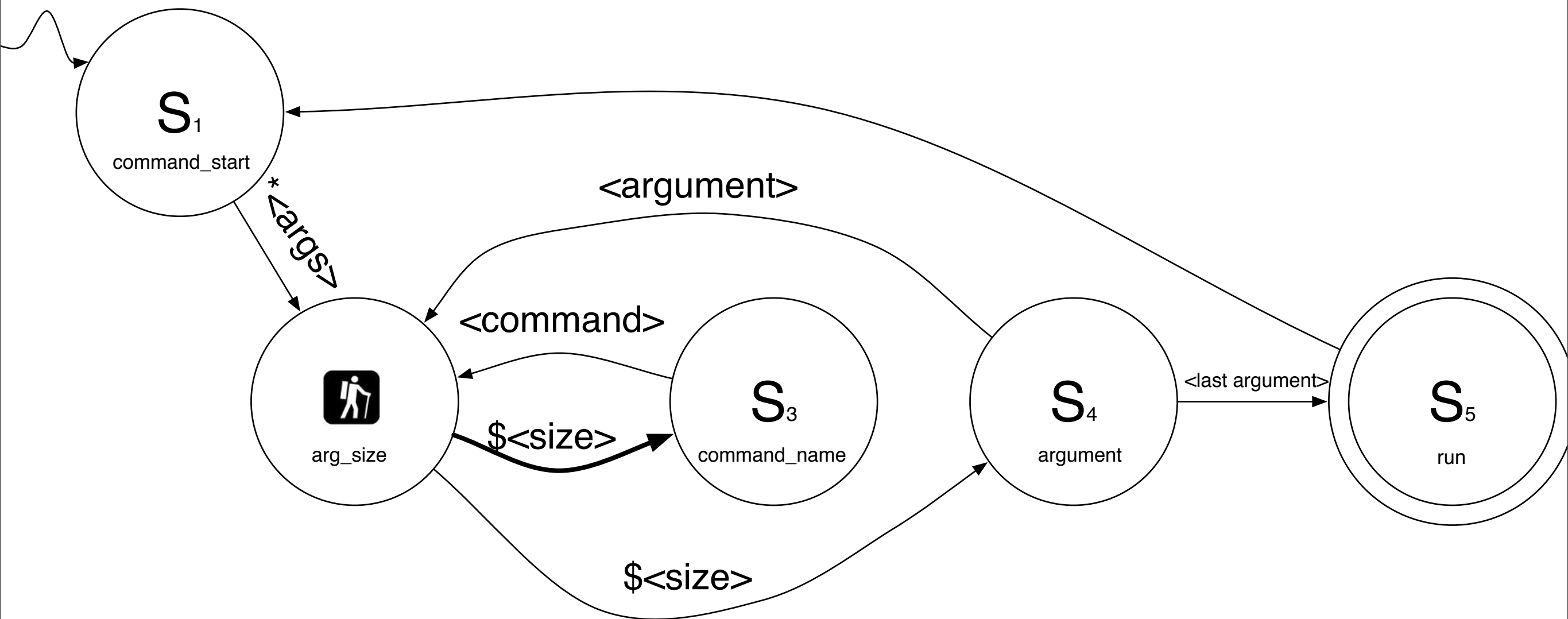
Redis Protocol



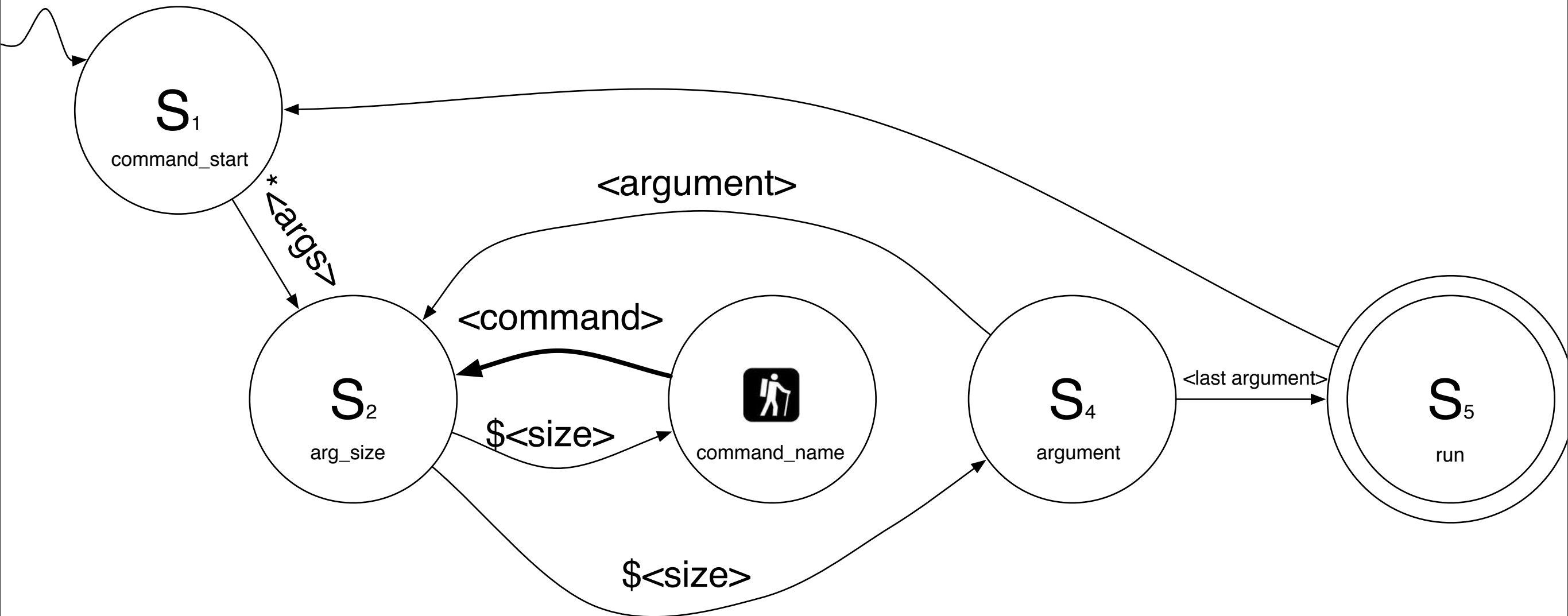

```
*DBG* <0.104.0> got {tcp,#Port<0.3714>,<<"*3\r\n">>} in state command_start
*DBG* <0.104.0> switched to state arg_size
```



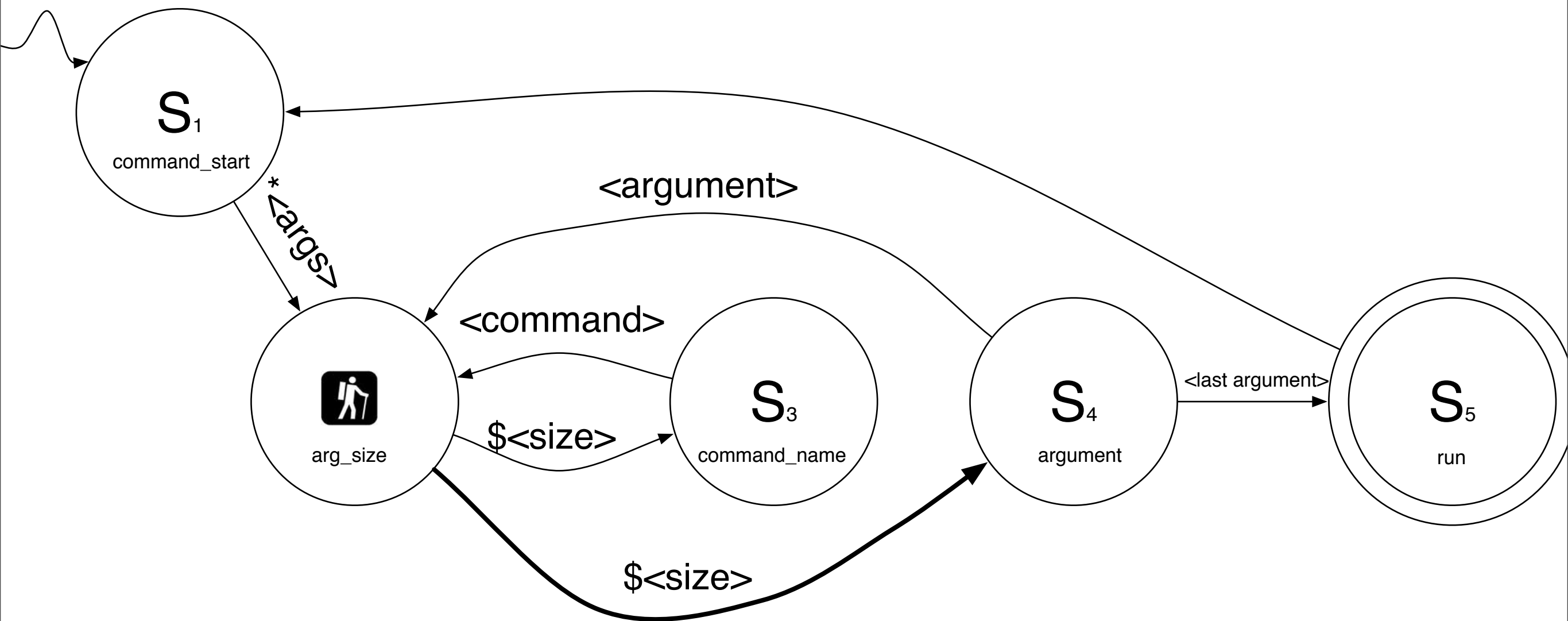
```
*DBG* <0.104.0> got {tcp,#Port<0.3714>,<<"$3\r\n">>} in state arg_size
*DBG* <0.104.0> switched to state command_name
```



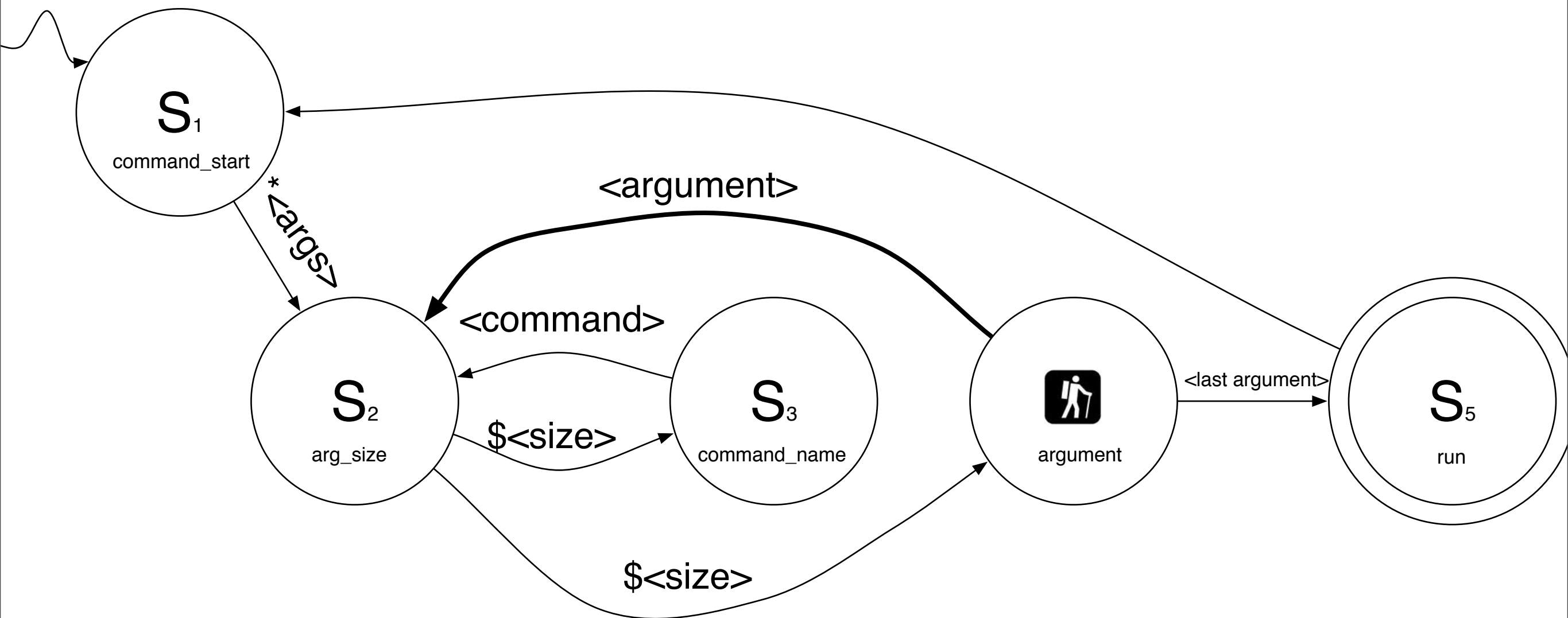
DBG <0.104.0> got {tcp,#Port<0.3714>,<<"set\r\n">>} in state command_name
DBG <0.104.0> switched to state arg_size



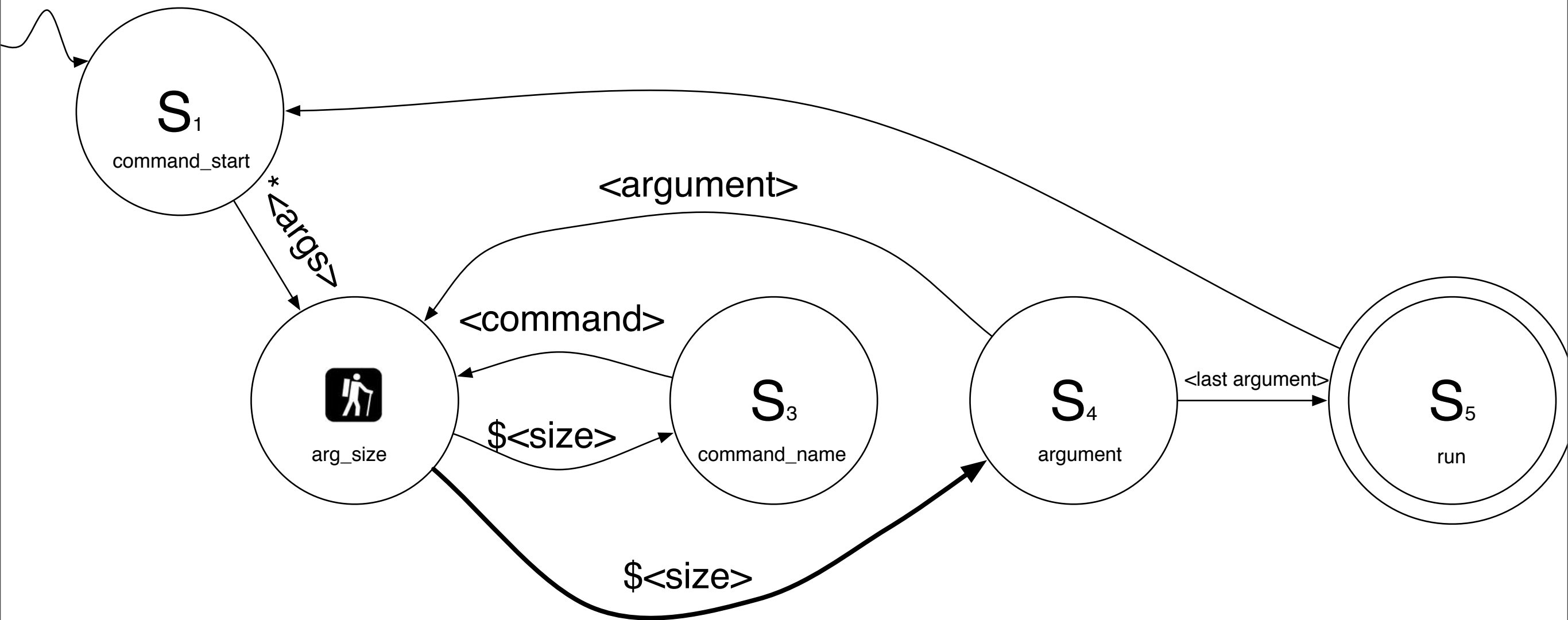
DBG <0.104.0> got {tcp,#Port<0.3714>,<<"\$4\r\n">>} in state arg_size
DBG <0.104.0> switched to state argument



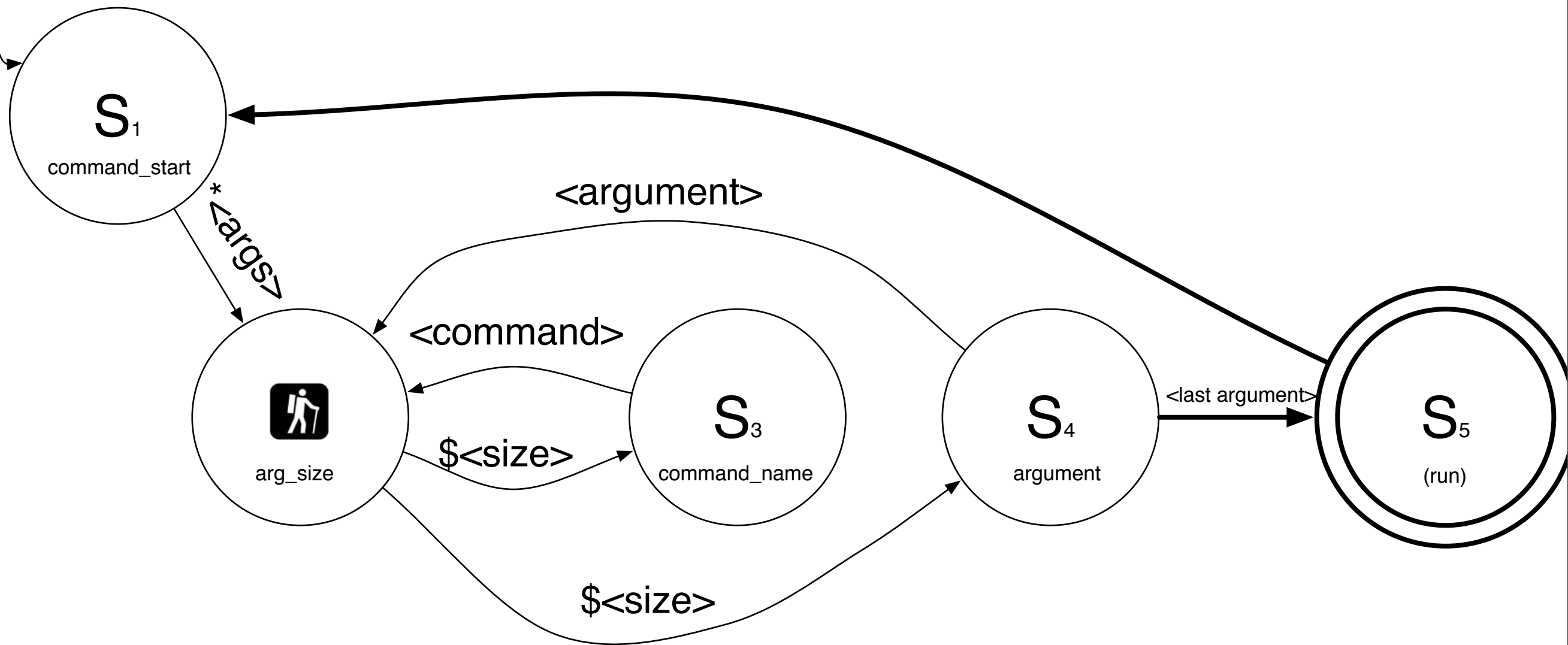
DBG <0.104.0> got {tcp,#Port<0.3714>,<<"lang\r\n">>} in state argument
DBG <0.104.0> switched to state arg_size



DBG <0.104.0> got {tcp,#Port<0.3714>,<<"\$5\r\n">>} in state arg_size
DBG <0.104.0> switched to state argument



DBG <0.104.0> got {tcp,#Port<0.3714>,<<"erlang\r\n">>} in state argument
DBG <0.104.0> switched to state command_start

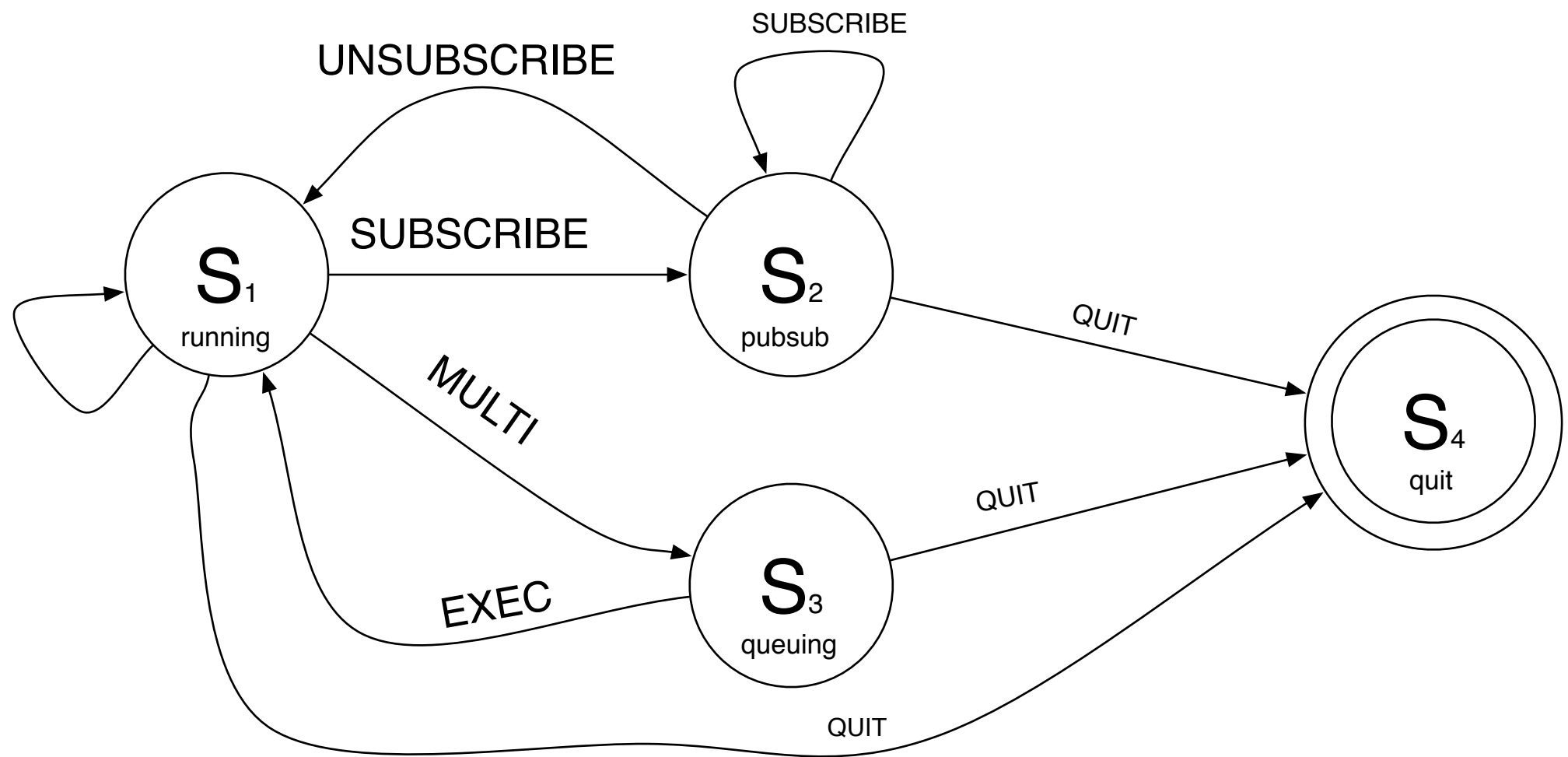


```
*DBG* <0.105.0> got cast {run,<<"SET">>,[<<"lang">>,<<"erlang">>]}
*DBG* <0.105.0> new state {state,#Port<0.3714>,  
                           'edis-db-0',0,56068,true,undefined,[],undefined}
```


edis_client.erl

```
120     edis_command_runner:run(State#state.command_runner,  
121         edis_util:upper(Command), []),  
122     {next_state, command_start, State, hibernate}
```

State Machine – Runner



command_runner

```
87 handle_cast({run, Cmd, Args}, State) ->
88     try
89         OriginalCommand = #edis_command{cmd = Cmd,
90                                         db = State#state.db_index,
91                                         args = Args},
92
93         Command = parse_command(OriginalCommand),
94
95         ok = edis_db_monitor:notify(OriginalCommand),
96
97         case {State#state.multi_queue, State#state.subscriptions} of
98             {undefined, undefined} -> run(Command, State);
99             {undefined, _InPubSub} -> pubsub(Command, State);
100             {_InMulti, undefined} -> queue(Command, State);
101             {_InMulti, _InPubSub} -> throw(invalid_context)
102         end
103     catch
```

command_runner

```
87 handle_cast({run, Cmd, Args}, State) ->
88     try
89         OriginalCommand = #edis_command{cmd = Cmd,
90                                         db = State#state.db_index,
91                                         args = Args},
92
93         Command = parse_command(OriginalCommand),
94
95         ok = edis_db_monitor:notify(OriginalCommand),
96
97         case {State#state.multi_queue, State#state.subscriptions} of
98             {undefined, undefined} -> run(Command, State);
99             {undefined, _InPubSub} -> pubsub(Command, State);
100             {_InMulti, undefined} -> queue(Command, State);
101             {_InMulti, _InPubSub} -> throw(invalid_context)
102         end
103     catch
```

command_runner

```
87 handle_cast({run, Cmd, Args}, State) ->
88     try
89         OriginalCommand = #edis_command{cmd = Cmd,
90                                         db = State#state.db_index,
91                                         args = Args},
92
93         Command = parse_command(OriginalCommand),
94
95         ok = edis_db_monitor:notify(OriginalCommand),
96
97         case {State#state.multi_queue, State#state.subscriptions} of
98             {undefined, undefined} -> run(Command, State);
99             {undefined, _InPubSub} -> pubsub(Command, State);
100             {_InMulti, undefined} -> queue(Command, State);
101             {_InMulti, _InPubSub} -> throw(invalid_context)
102         end
103     catch
```

edis_db

```
67 run(Db, Command, Timeout) ->
68   try gen_server:call(Db, Command, Timeout) of
69     ok -> ok;
70     {ok, Reply} -> Reply;
71     {error, Error} ->
72       throw(Error)
73   catch
74     _:{timeout, _} ->
75       throw(timeout)
76   end.
```

edis_db

```
67 run(Db, Command, Timeout) ->
68   try gen_server:call(Db, Command, Timeout) of
69     ok -> ok;
70     {ok, Reply} -> Reply;
71     {error, Error} ->
72       throw(Error)
73   catch
74     _:{timeout, _} ->
75       throw(timeout)
76   end.
```

edis_db

```
214 handle_call(#edis_command{cmd = <<"MSET">>, args = KVs},
215             _From, State) ->
216     Reply =
217         (State#state.backend_mod):write(
218             State#state.backend_ref,
219             [{put, Key,
220              #edis_item{key = Key, encoding = raw,
221                       type = string, value = Value}}
              || {Key, Value} <- KVs]),
222     {reply, Reply, stamp([K || {K, _} <- KVs], write, State)};
```


edis_db

```
214 handle_call(#edis_command{cmd = <<"MSET">>, args = KVs},
215             _From, State) ->
216     Reply =
217     (State#state.backend_mod):write(
218     State#state.backend_ref,
219     [{put, Key,
220     #edis_item{key = Key, encoding = raw,
221     type = string, value = Value}}
222     || {Key, Value} <- KVs]),
222     {reply, Reply, stamp([K || {K, _} <- KVs], write, State)};
```

```
35 write(#ref{db = Db}, Actions) ->
36     ParseAction = fun({put, Key, Item}) ->
37         {put, Key, erlang:term_to_binary(Item)};
38         (Action) -> Action
39     end,
40     eleveldb:write(Db, lists:map(ParseAction, Actions), []).
41
```

```
35 write(#ref{db = Db}, Actions) ->
36     ParseAction = fun({put, Key, Item}) ->
37         {put, Key, erlang:term_to_binary(Item)};
38         (Action) -> Action
39     end,
40     eleveldb:write(Db, lists:map(ParseAction, Actions), []).
41
```

edis_command_runner.erl

```
700 run(C = #edis_command{result_type = ResType,  
701      timeout = Timeout, hooks = Hooks}, State) ->  
702      Res = edis_db:run(State#state.db, C);
```

edis_command_runner.erl

```
716     case ResType of
717         ok -> tcp_ok(State);
718         string -> tcp_string(Res, State);
719         bulk -> tcp_bulk(Res, State);
720         multi_bulk -> tcp_multi_bulk(Res, State);
721         number -> tcp_number(Res, State);
722         boolean -> tcp_boolean(Res, State);
723         float -> tcp_float(Res, State);
724         sort -> tcp_sort(Res, State);
725         zrange ->
726             [_Key, _Min, _Max, ShowScores, Limit] = C#edis_command.args,
727             tcp_zrange(Res, ShowScores, Limit, State)
728     end.
```

LevelDB

- Stores arbitrary byte arrays
- Data is stored sorted by key
- Three operations: Put/Get/Delete
- Multiple changes in atomic batch operations
- Data is automatically compressed

(LevelDB is not the only backend, currently, in-memory and HanoiDB)

So what is different?

- SAVE, BGSAVE, LASTSAVE: database dependent
- MULTI doesn't support cross-db/non-db
- SLAVEOF not fully supported
- Encoding is inefficient

Master/Slave Replication

Multi-Master Behavior

Replication



Slave sends SYNC
command



Master (may) flush RDB to
disk



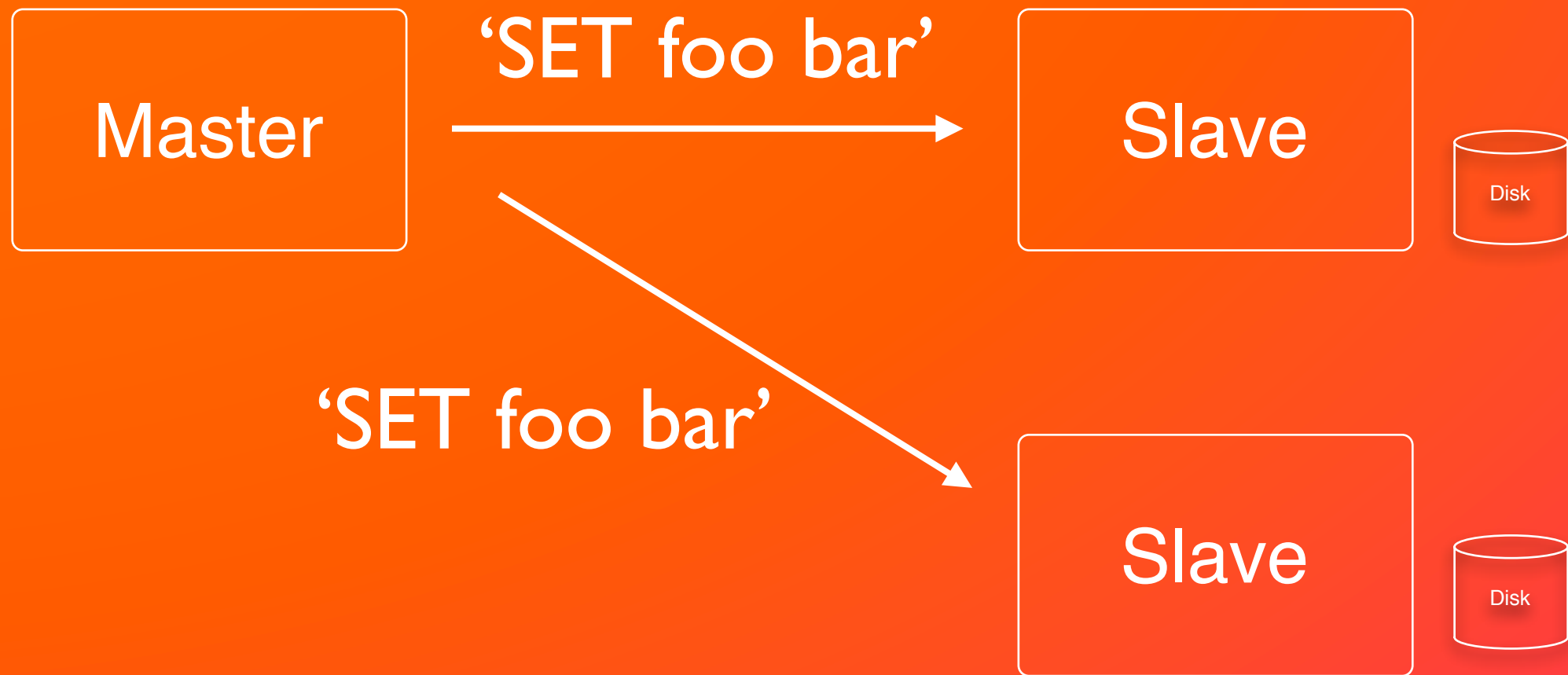
**Master sends database file
to slave**



Slave saves file as its
new .rdb file



Slave now **MONITORs**
all commands



Masters can have multiple slaves



Slaves can have slaves

‘SET foo bar’



...‘SET foo bar2’



Slaves can accept writes
(this can be disabled)

‘SET foo1 s1’



Master



... ‘SET foo2 s2’



Master

‘GET foo2’ -> s2

‘GET foo1’ -> s1

No multi-master

‘SET foo1 s1’

... **‘SET foo2 s2’**



‘GET foo2’ -> s2

‘GET foo1’ -> s1

No multi-master

'SET foo1 s1'

... 'SET foo2 s2'

Master

Master



'GET foo2' -> s2

'GET foo1' -> s1

No multi-master

'SET foo1 s1'

... 'SET foo2 s2'



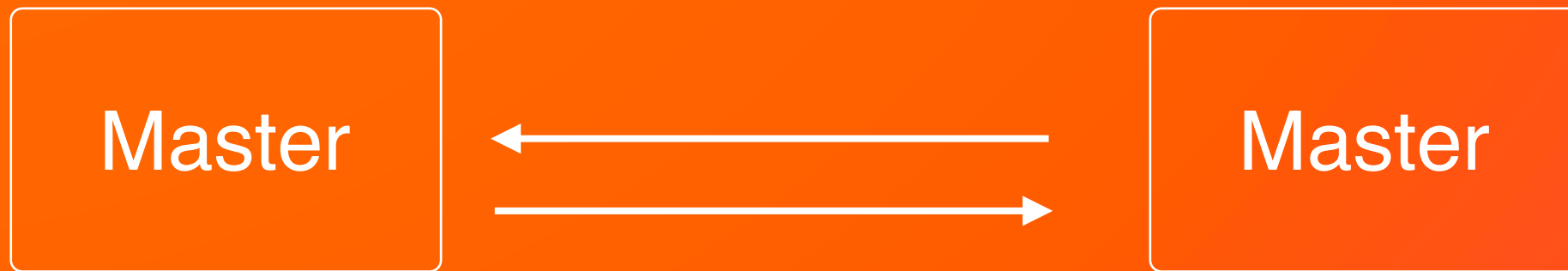
'GET foo2' -> s2

'GET foo1' -> s1

No multi-master in Redis (yet)

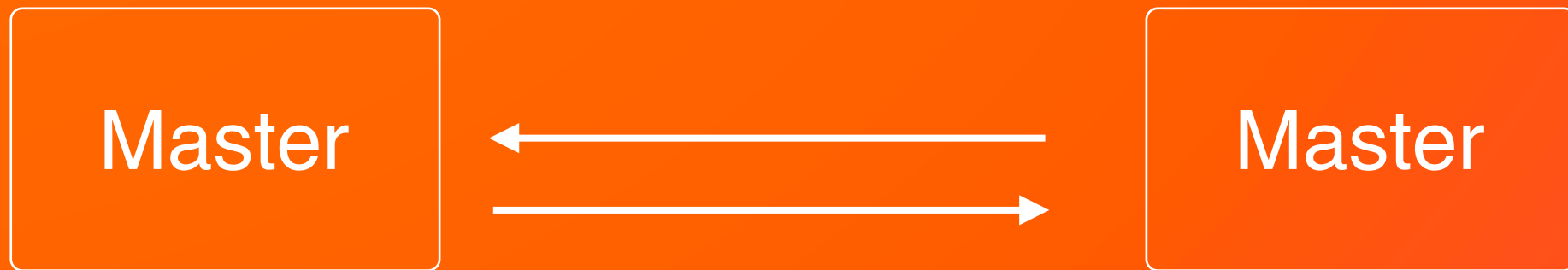
So how does Edis Multi-
master work?

‘SET foo l sl’



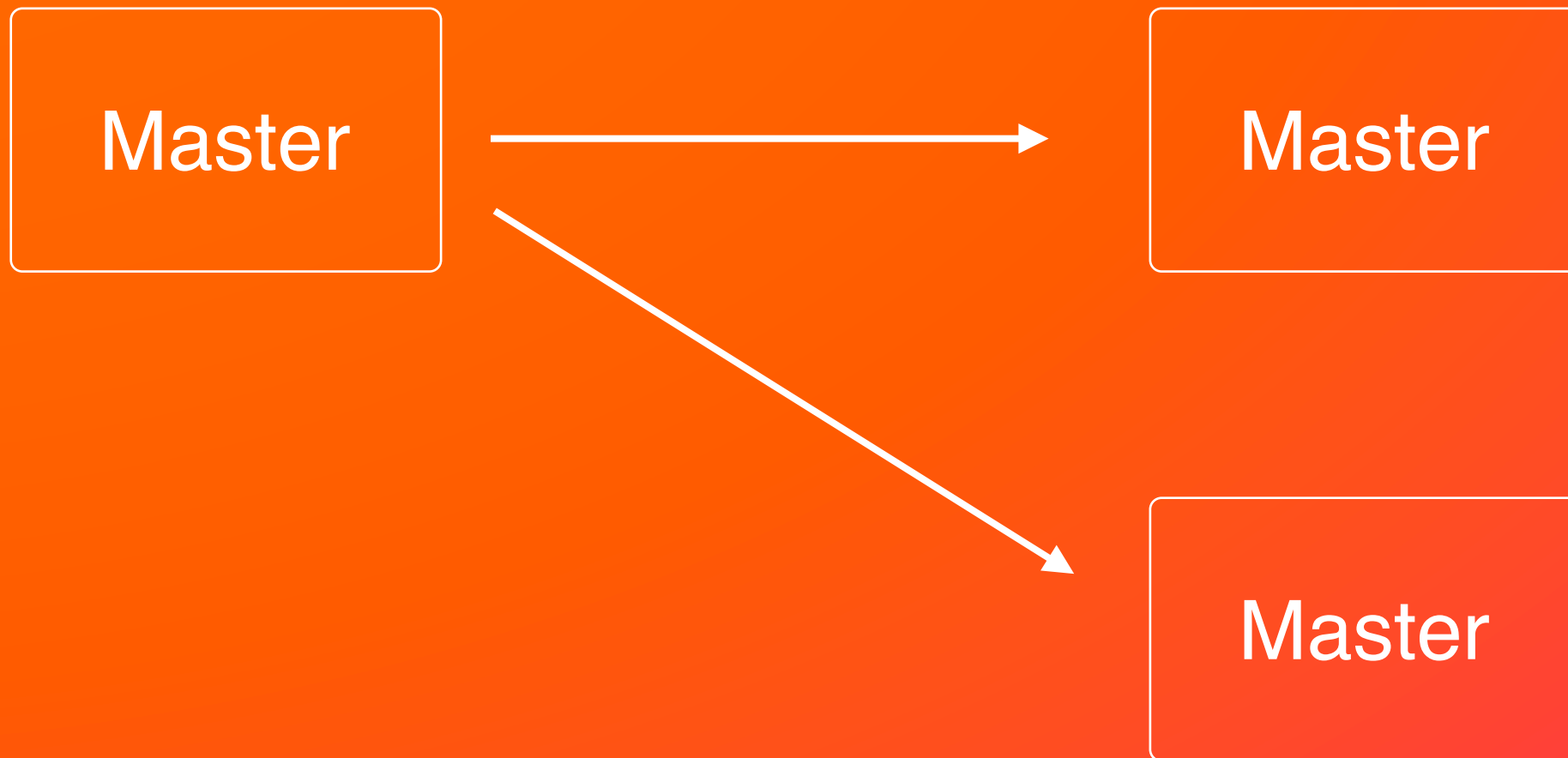
Uses `gen_server:abcast/2`

‘SET foo 1 s1’



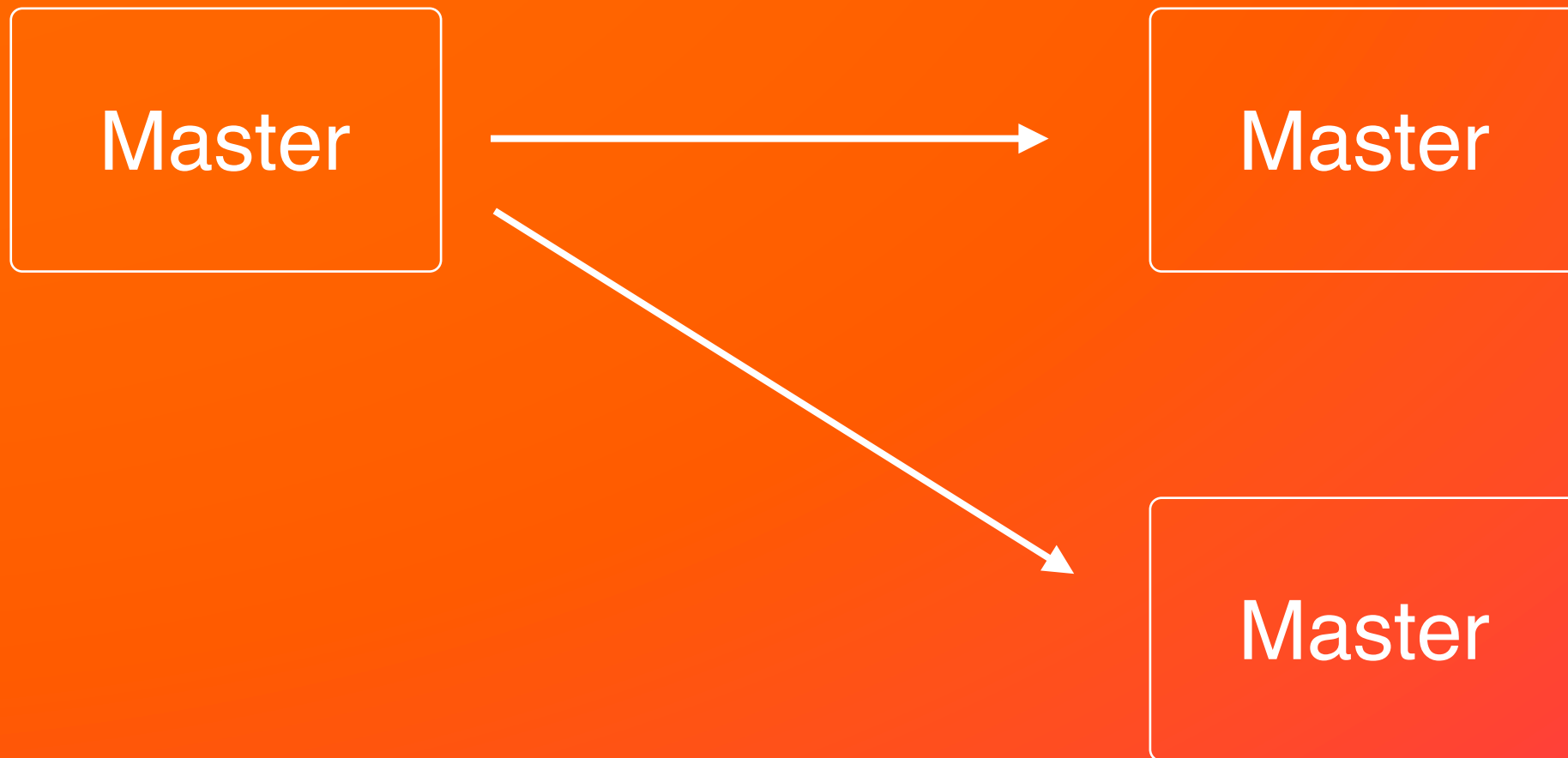
Broadcasts to all nodes

‘SET foo | sl’



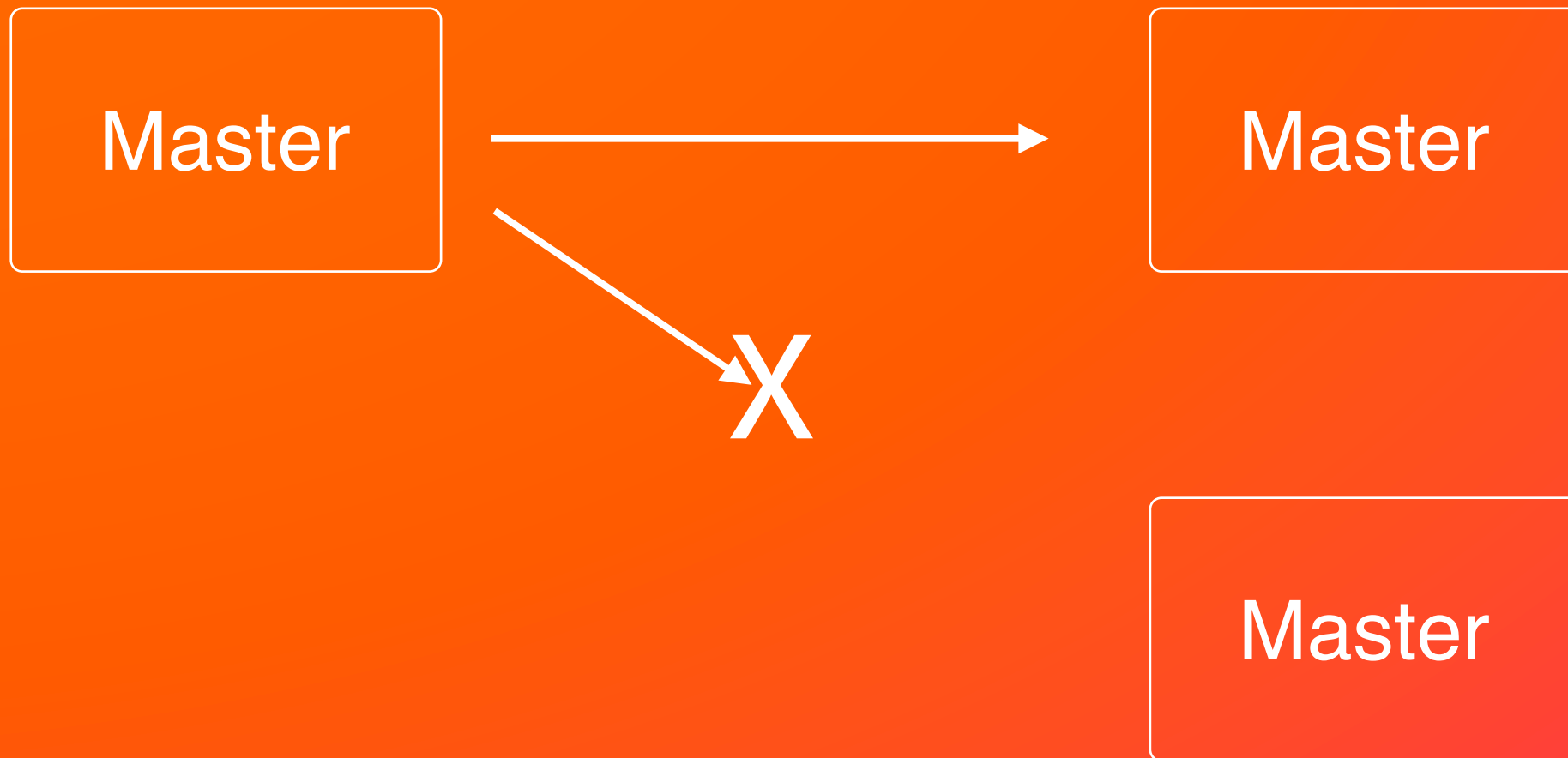
Broadcasts to all nodes

‘SET foo | sl’



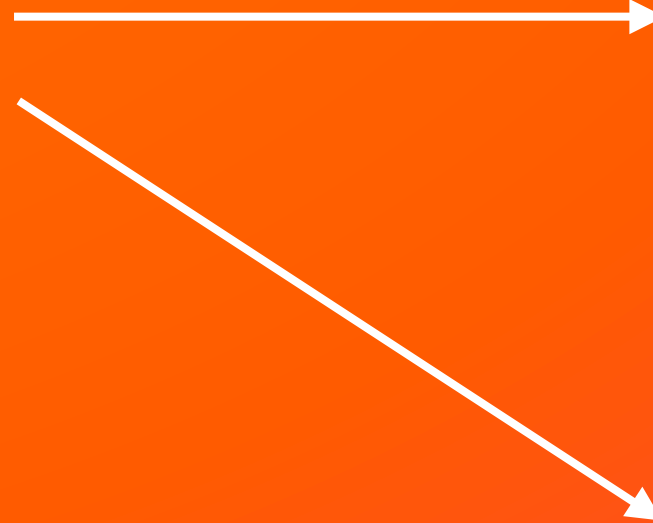
Uses vector clocks to
resolve key versions

‘SET foo 1 sl’



Does not update
disconnected nodes on
reconnect

'SET fool sl'



'fool -> sl'



fool -> []

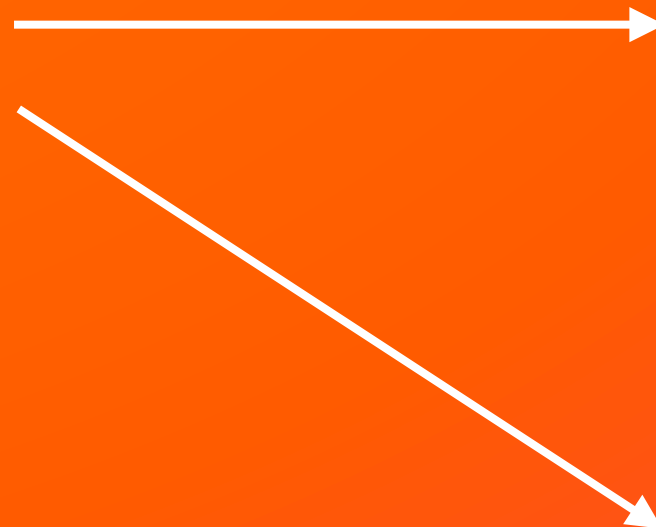


Does not update
disconnected nodes on
reconnect

‘SET foo l sl’



‘foo l -> sl’



foo l -> []



However disconnected nodes will
detect that keys are out of date
since last connection

What is Edis's CAP?

Consistent

Available

Partition-Tolerant

What is Edis's CAP?

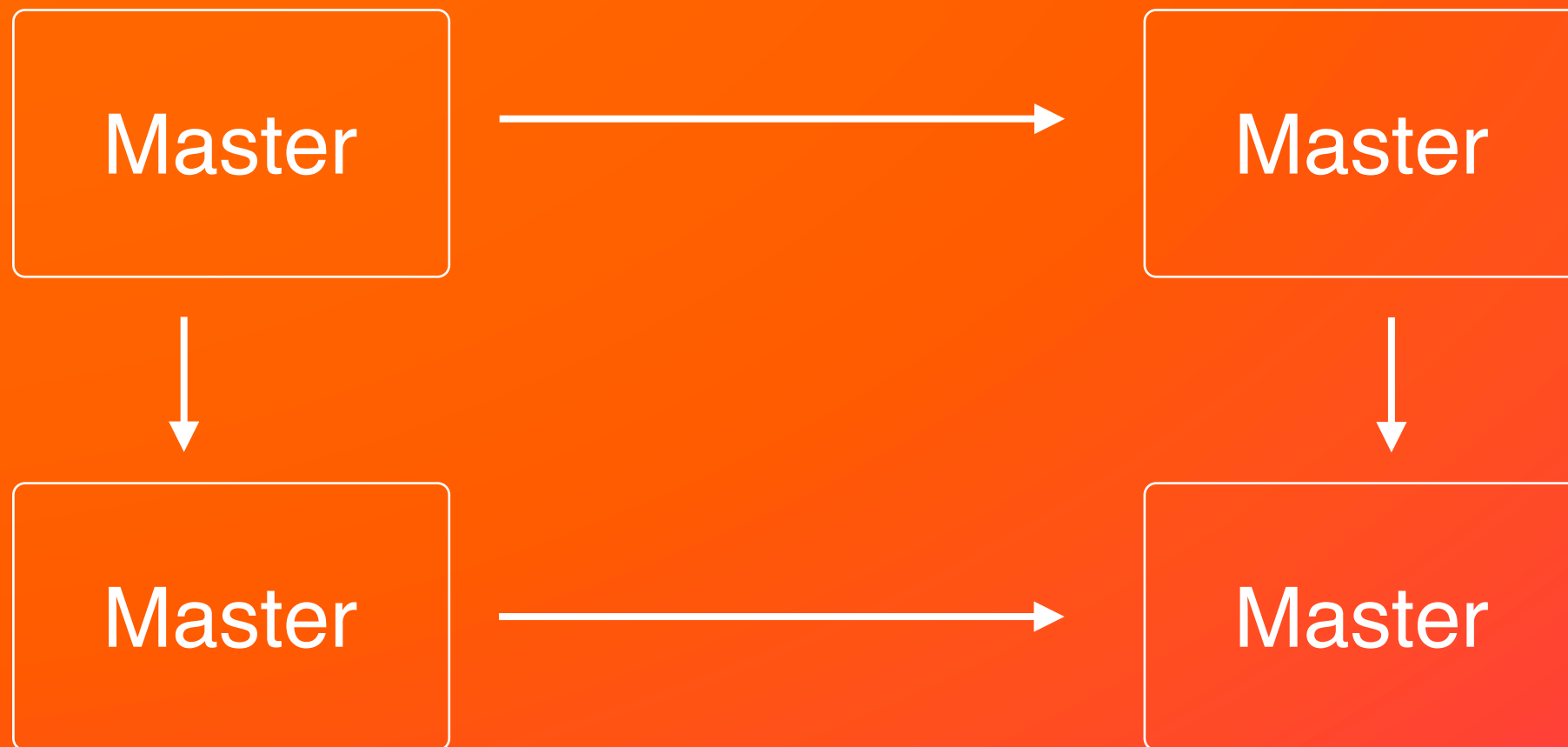
Consistent

Available

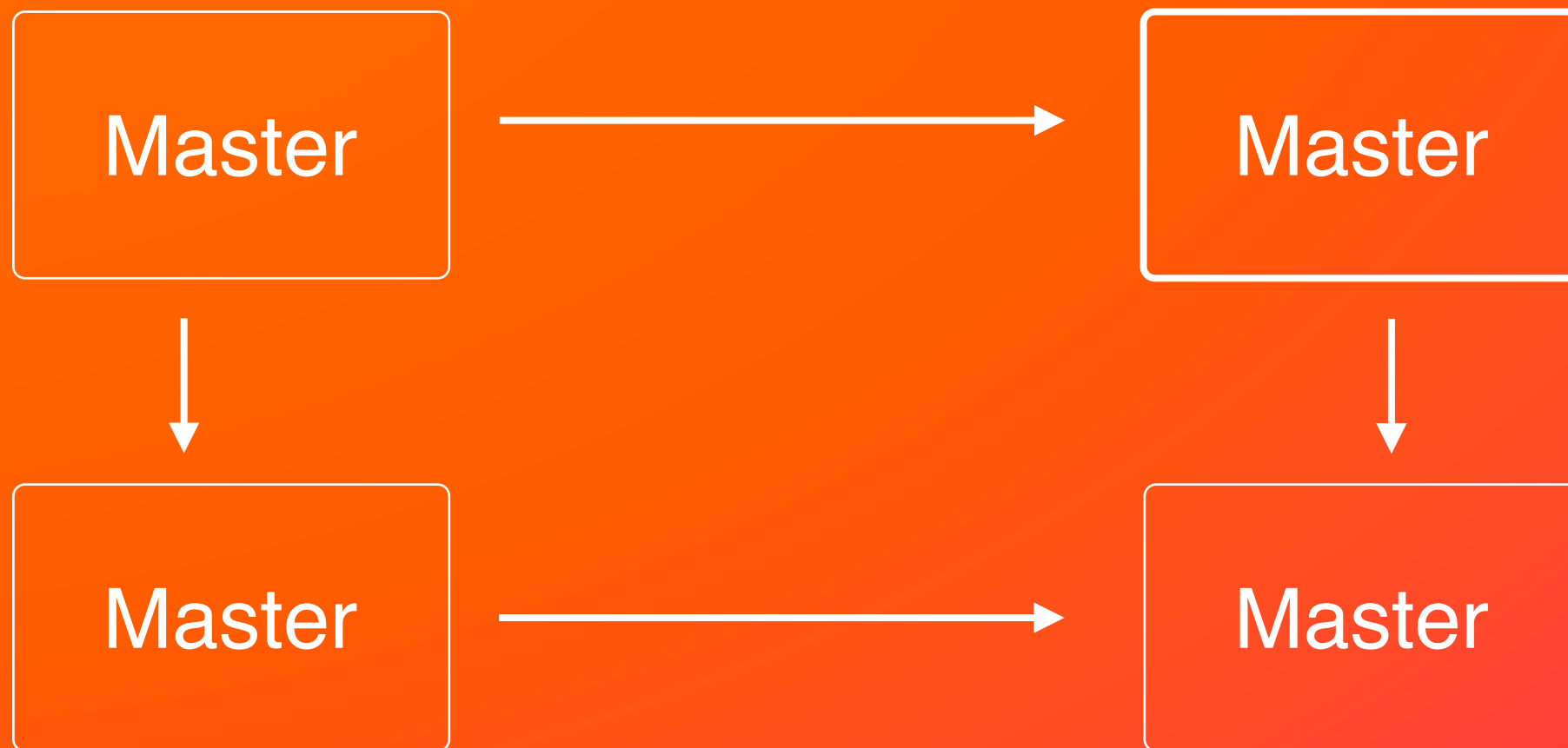
Partition-Tolerant

What's next?

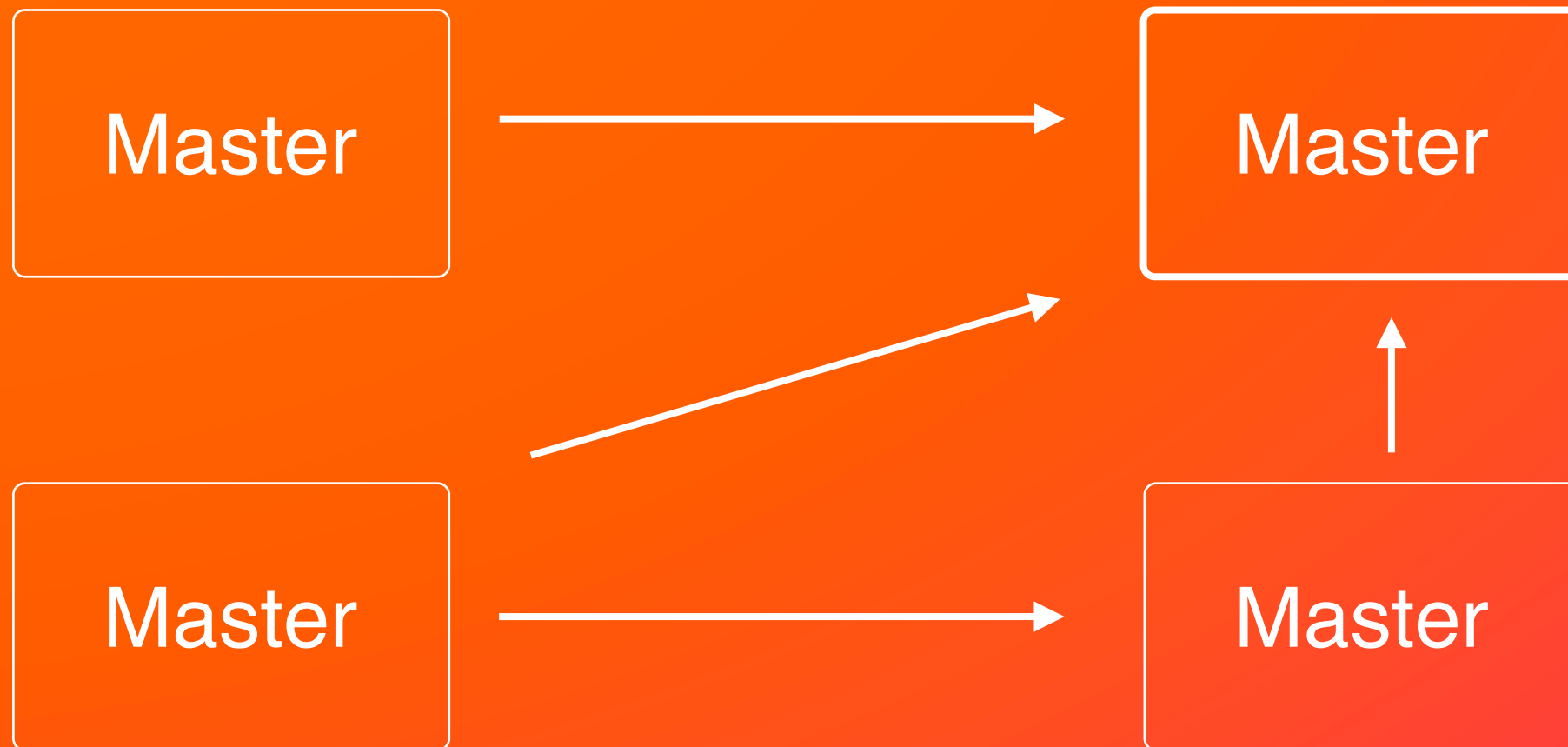
Computation Nodes



Leader Election using
gen_leader



Leader Election using
gen_leader



now INCR, POP, etc occur
on leader only

Replication on reconnect

- Master asks leader for SYNC but disconnects after database received
- New node can either block before replying or background replicate (Currently replication is background)

- Finish Replication for all datatypes
- Post commit hooks
- Redis cluster implementation

How can you help?

- Look at the Edis source
- Lots of Benchmarks w/ common test
- Experiment with it

<https://github.com/inaka/edis>

Thanks to Joachim Nilsson
for his Uppsala University thesis work
applied to Edis while at Inaka

Thank you / Questions?

<http://inaka.net>

@chaddepue

chad@inaka.net || chad@whisper.sh

inaka

Other slides

Performance (Operations/second)

	Redis	In-Memory Edis	% slower
PING (inline)	120,734	40,741	296%
PING	129,892	32,956	394%
MSET (10 keys)	73,825	6,662	1,108%
SET	135,160	22,051	613%
GET	134,282	23,127	581%
INCR	138,916	24,421	569%
LPUSH	137,990	21,397	645%
LPOP	130,769	22,728	575%
SADD	135,160	21,860	618%
SPOP	132,456	25,707	515%
LRANGE (first 100 elements)	65,362	1,783	3,667%

Performance (Operations/second)

	Redis	LevelDB Edis	% slower
PING (inline)	120,734	41,152	293%
PING	129,892	32,419	401%
MSET (10 keys)	73,825	6,058	1,219%
SET	135,160	20,726	652%
GET	134,282	21,463	626%
INCR	138,916	17,930	775%
LPUSH	137,990	226	61,105%
LPOP	130,769	229	57,092%
SADD	135,160	9,003	1,501%
SPOP	132,456	1,298	10,205%
LRANGE (first 100 elements)	65,362	644	10,143%