

# The Design of Scalable Distributed (SD) Erlang

*Natalia Chechina, Amir Ghaffari, Phil Trinder,  
and RELEASE team*

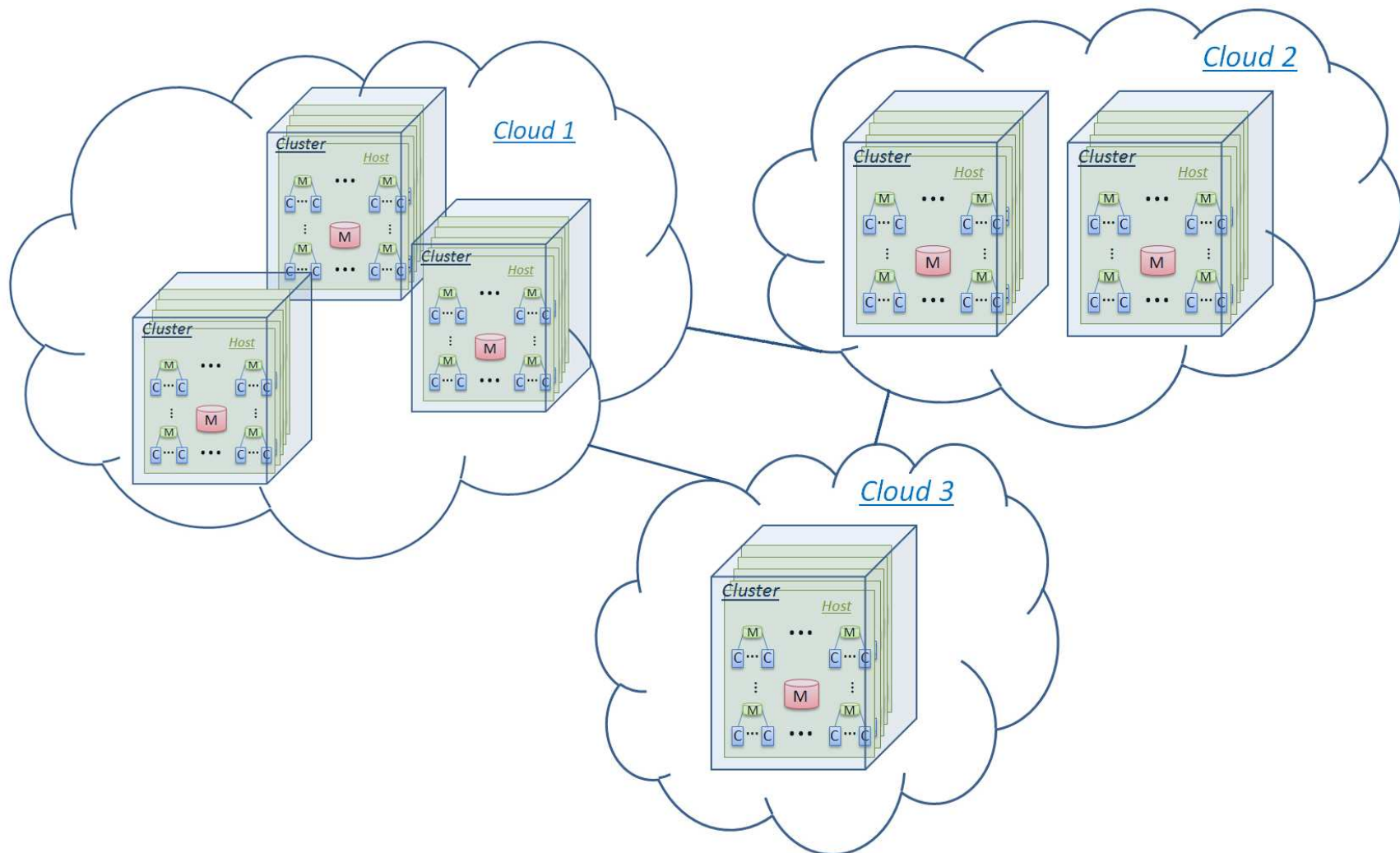
June 15, 2012



# Design Approach

- Typical hardware architecture
- Scaling
  - Persistent data structures
  - In-memory data structures
  - Computation
- Security

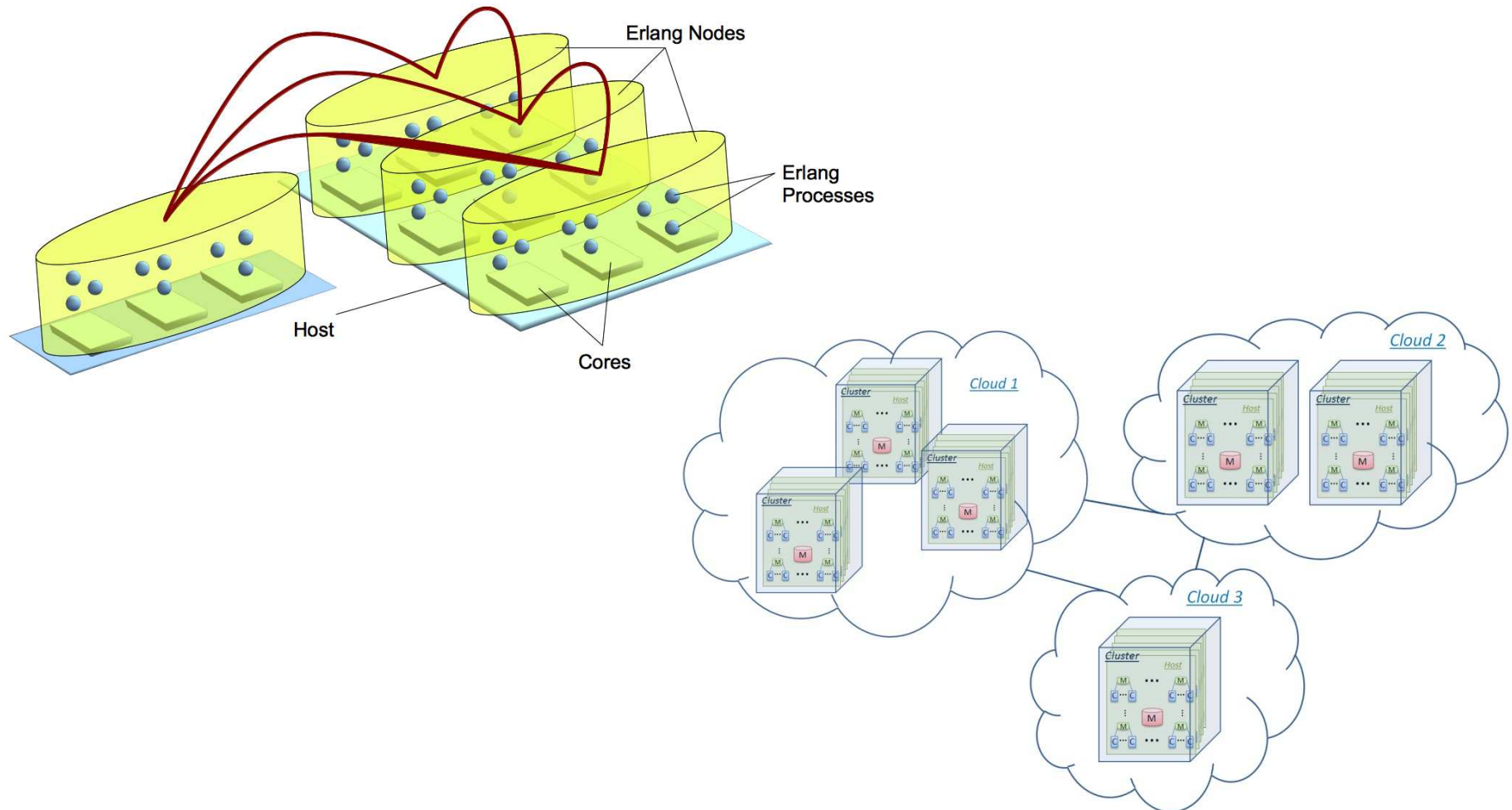
# Typical architecture – $10^5$ cores



# Design Approach

- Typical hardware architecture
- Scaling
  - Persistent data structures
  - In-memory data structures
  - Computation
- Security

# Distributed Erlang



# Scaling Computation

- Network Scalability
  - All to all connections are not scalable onto 1000s of nodes
  - Aim: Reduce connectivity
- Semi-explicit Placement
  - Becomes not feasible for a programmer to be aware of all nodes
  - Aim: Automatic process placement in groups of nodes

# Design Principles

## General:

- Working at Erlang level as far as possible
- Preserving the Erlang philosophy and programming idioms
- Minimal design changes

# Design Principles

## Reliable Scalability:

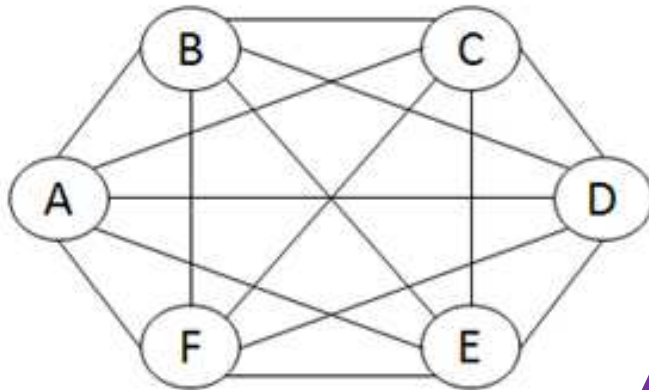
- Avoiding global sharing
- Avoiding explicit prescription
- Introducing an abstract notion of communication architecture
- Keeping Erlang reliability model unchanged as far as possible



# Network Scalability

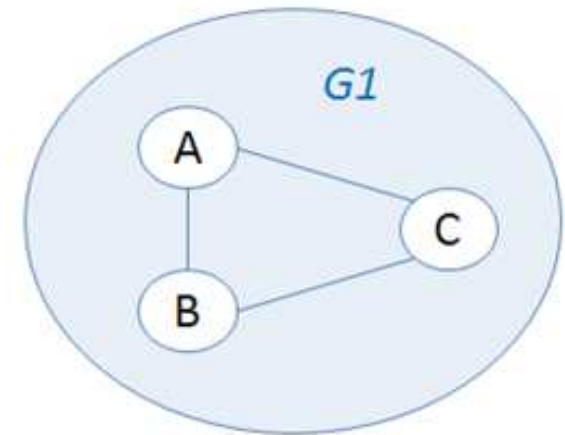
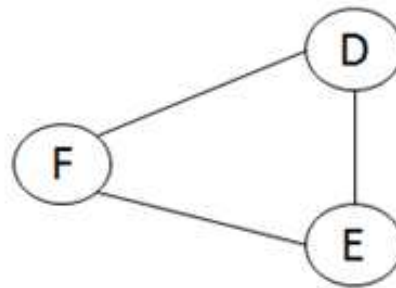
- Grouping nodes in Scalable groups (s\_groups)
  - transitive connections with nodes of the same s\_group
  - non-transitive connections with other nodes
- Types of s\_groups:
  - Hierarchical
  - Overlapping
  - Partition
- Using s\_group variables instead of global variables: *Var@Group*

# Creating an s\_group



A: `new_s_group(G1, [A, B, C]).`

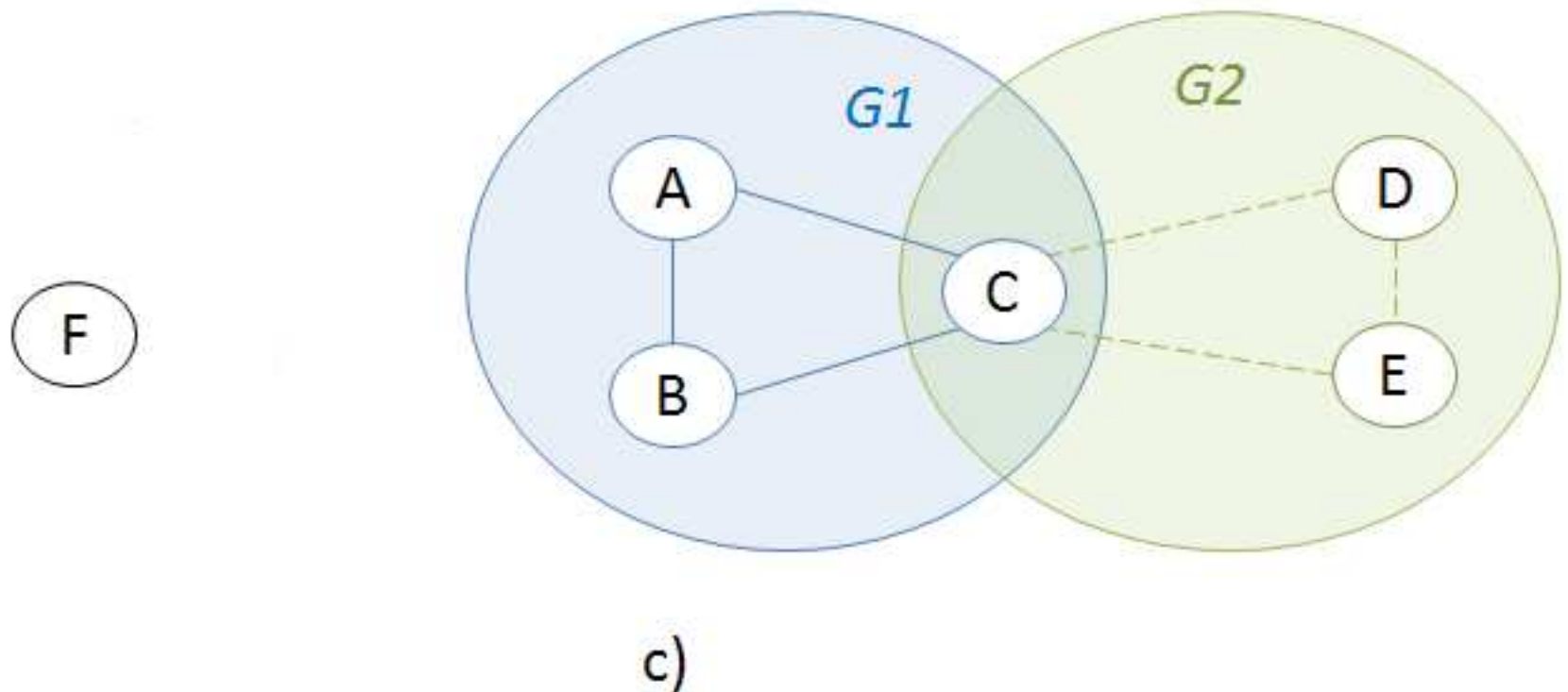
a)



b)

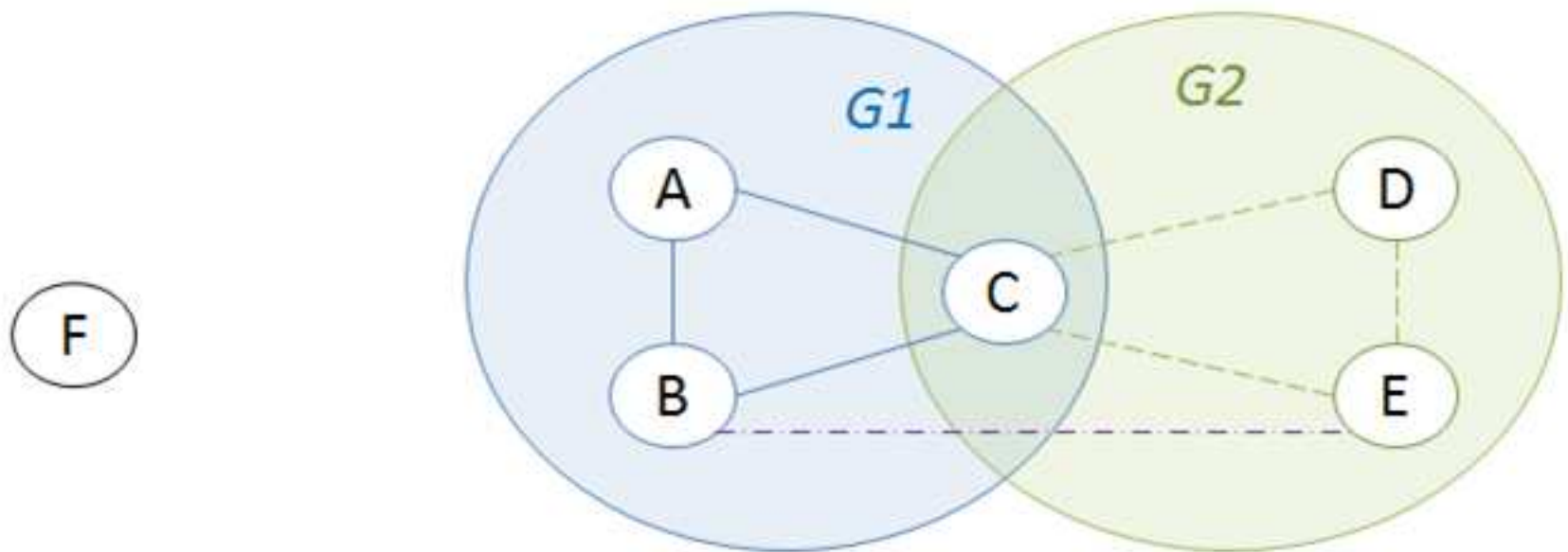
# Overlapping Groups & Non-transitive Connections

C: new\_s\_group(G2, [C, D, E]).



# Any to Any Connection

B: spawn(E, f).



d)

# s\_group Functions (1)

- Creating a new s\_group

`new_s_group($GroupName, [Node]) -> true | {error, ErrorMsg}`

- Deleting an s\_group

`del_s_group($GroupName) -> true | {error, ErrorMsg}`

- Adding new nodes to an existing s\_group

`add_node_s_group($GroupName, [Node]) -> true | {error, ErrorMsg}`

- Removing nodes from an existing s\_group

`remove_node_s_group($GroupName, [Node]) -> true | {error, ErrorMsg}`

# s\_group Functions (2)

- Monitoring all nodes of an s\_group

`monitor_s_group(S_GroupName) -> ok | {error, ErrorMessage}`

- Sending a message to all nodes of an s\_group

`send_s_group(S_GroupName, Msg) -> Pid | {badarg, Msg} {error, ErrorMessage}`

- Listing nodes of a particular s\_group

`s_group_nodes(S_GroupName) -> [Node] | {error, ErrorMessage}`

- Listing s\_groups that a particular node belongs to

`node_s_group_info(Node) -> [S_GroupName]`

# Scaling Computation

- Semi-explicit Placement
  - Becomes not feasible for a programmer to be aware of all nodes and place each of them explicitly
  - Aim: Automatic process placement

# chose\_node/1

chose\_node(Restrictions) -> node()

Restrictions = [Restriction]

Restriction = {s\_group, S\_Group}

| {min\_dist, MinDist :: integer() >= 0}

| {max\_dist, MaxDist :: integer() >= 0}

| {ideal\_dist, IdealDist :: integer() >= 0}

start() ->

TargetNode = chose\_node({s\_group, S\_Group},  
{ideal\_dist, IdealDist}),

spawn(TargetNode, fun() -> loop() end).



# Thank you!