

# Quo vadis Linux-HA?

## Developments in Linux Clustersoftware

Michael Schwartzkopff  
misch@schwartzkopff.org

# Background



# Linux Clustersoftware

- Linux Virtual Server
  - Scales nicely...
  - ... but is not high available.
- Linux-HA
  - Offers high availability ...
  - ... but does not scale.
- Code is not maintained any more.
- Is being replaced by a collection of other programs.

# Linux-HA Version 1 (heartbeat)

- Two servers exchange *heartbeats*.
- If the standby server does not receive *heartbeats* from the active node any more it starts the services.
- Configuration in a plain text file *haresources*.
- Cons:
  - No monitoring of the services.
  - No management of the resources.

# Linux-HA Version 2 (CRM)

- Management of the resources by a *Cluster Resource Manager (CRM)*:
  - Aktive monitoring of the reosources by the cluster software.
  - Up to 16 nodes in a cluster.
  - Resources free moveable between the nodes.
  - Communication in the cluster via hear tbeat.
  - Constraints determine which resource should run on what node.

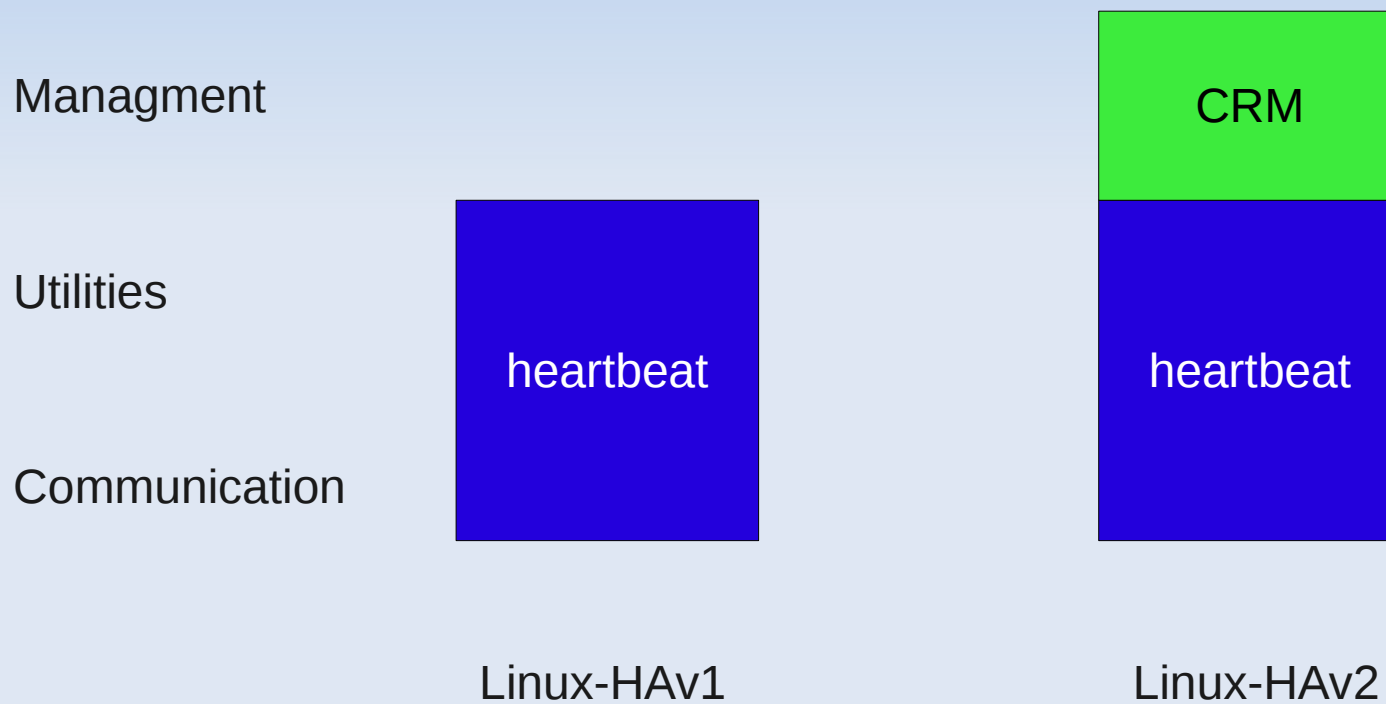
# Constraints in Linux-HAv2

- **Order**  
A resource should start before or after an other resource.
- **Colocation**  
A resource should run on the same node as an other resource.
- **Location**  
A resource should run on a node with specific attributes.

# Features of Linux-HAv2?

- Pros
  - Monitoring of the resources.
- Cons
  - Configuration in XML.
  - Administration via the command line.
  - GUI not really usable.
  - Not maintained any more after 2007 (version 2.1.4).

# Linux-HA

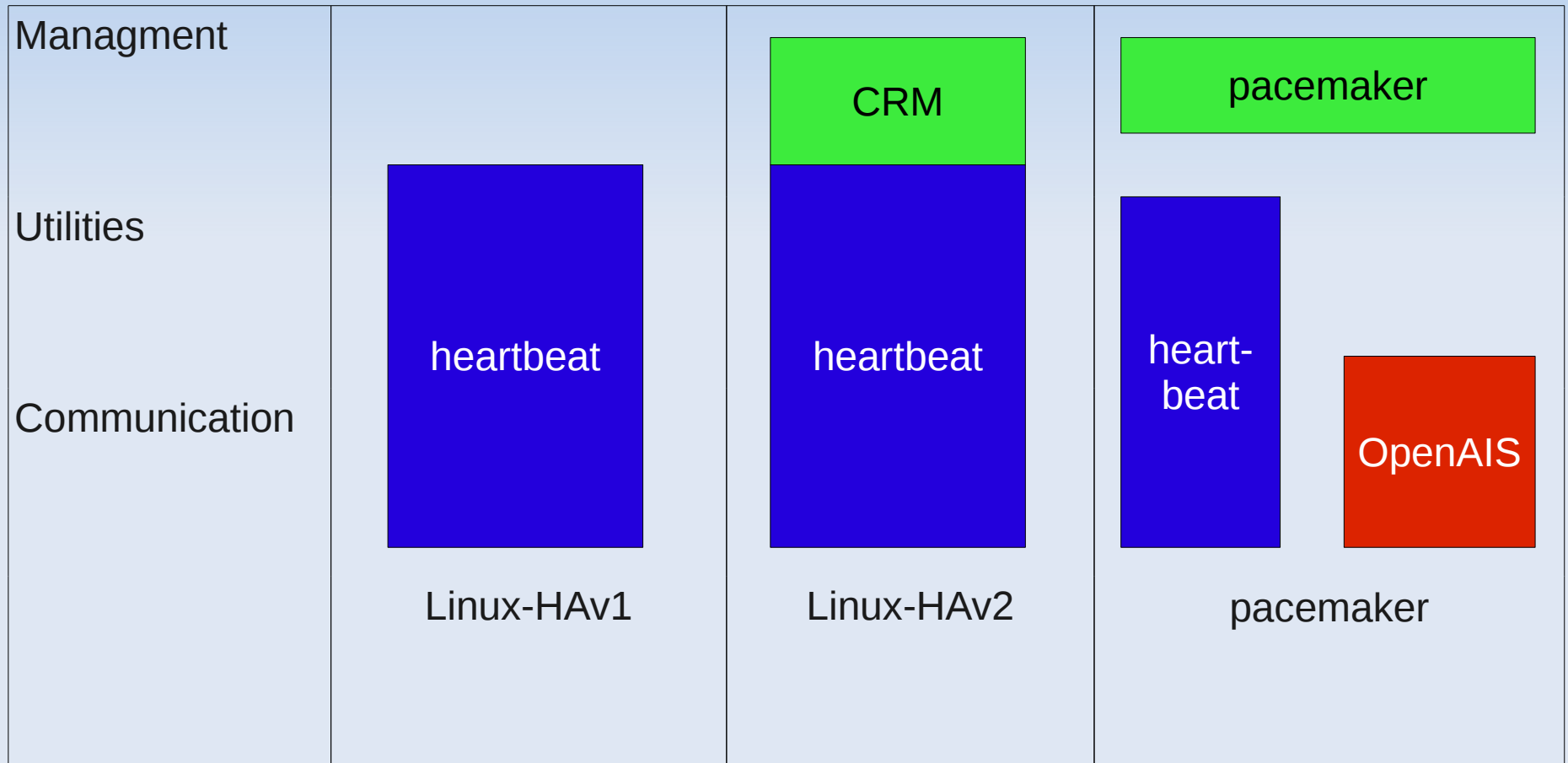




# Restart

- If the CRM uses heartbeat only for the communication in the cluster, this could be done also by another program.
- The other software was OpenAIS.
- The developers organized the CRM code in a separate project called pacemaker. It uses either heartbeat or OpenAIS.
- Some utilities from the heartbeat-package were still necessary.

# pacemaker



# pacemaker

- Uses heartbeat or OpenAIS.
- Suitable GUI.
- Own subshell additional to the CLI.
- Sandboxes for testing.
  - First test what would happen, if ...
    - only afterwards activate the new configuration.
- Exact history, what happened why.

# pacemaker: The GUI

The screenshot displays the Pacemaker GUI interface. The window title is "Pacemaker GUI <@opensuse1>". The menu bar includes "Connection", "View", "Shadow", "Tools", and "Help". The left sidebar shows a tree view under "Live" with categories: Configuration (CRM Config, Resource Defaults, Operation Defaults), Nodes, Resources, Constraints, and Management (selected).

The main area shows a table of cluster components:

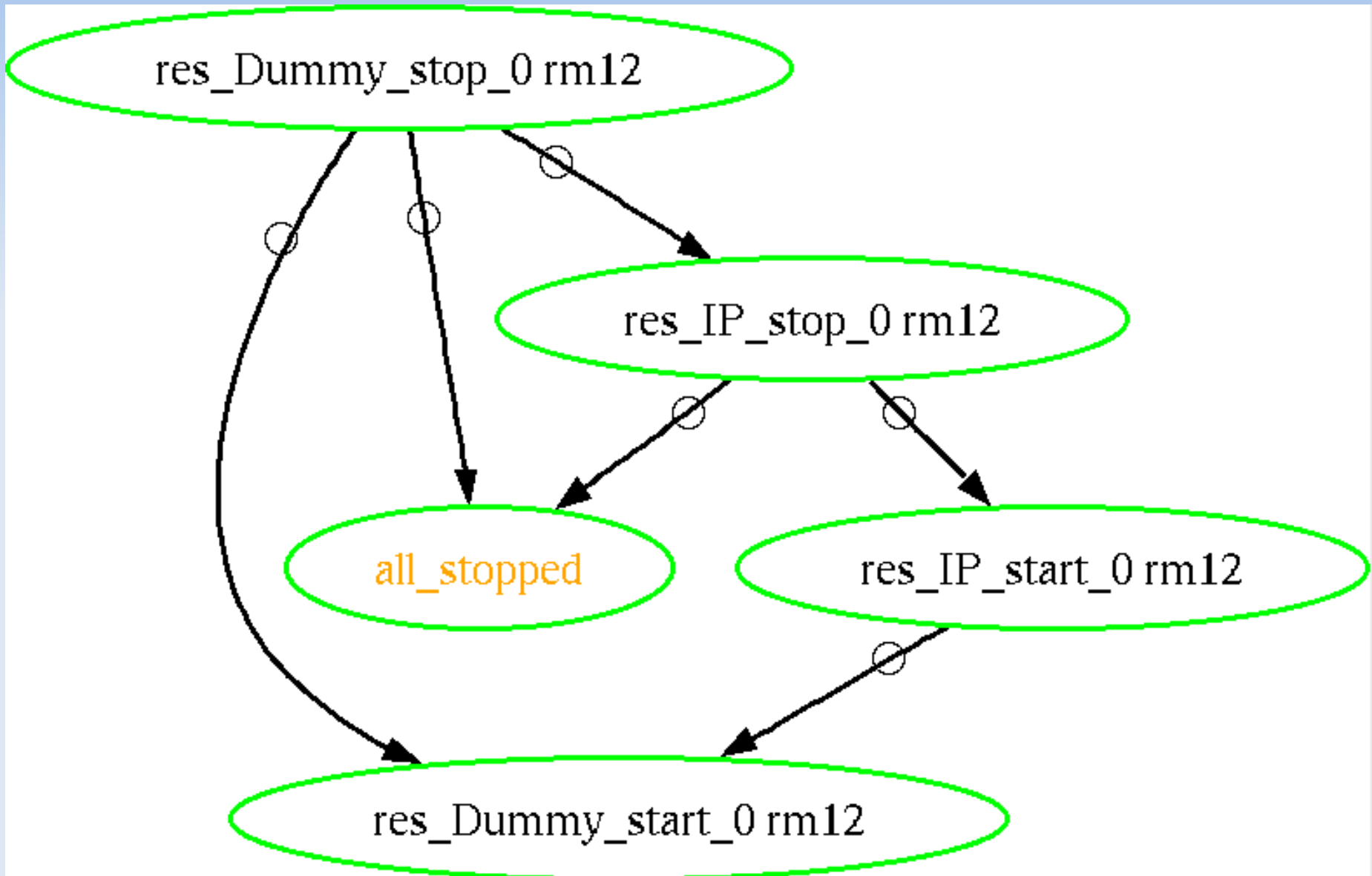
Name	Status	Details
Cluster	● have quorum	Openais & Pacemaker
opensuse1	● online (dc)	
opensuse2	● online	
Resources	●	
clonePingd	● clone	
resPingd:0	● running on ['opensuse1']	ocf::pacemaker:pingd
resPingd:1	● running on ['opensuse2']	ocf::pacemaker:pingd
groupWebserver	● group	
resIP	● running on ['opensuse1']	ocf::heartbeat:IPAddr2
resApache	● not running	ocf::heartbeat:apache

Below the table, the Migration Threshold is set to 1000000. A log table shows the following entries:

Call ID	Operation	Interval	Return Code	Status	Last Run	Exec Time	Queue Time	Last Return Code Change
21	probe		ok (rc=0)	complete	Tue Oct 13 16:31	180ms	10ms	Tue Oct 13 16:31:13 2009
23	stop		ok (rc=0)	complete	Tue Oct 13 16:31	320ms	0ms	Tue Oct 13 16:31:30 2009
24	start		ok (rc=0)	complete	Tue Oct 13 16:31	390ms	10ms	Tue Oct 13 16:31:33 2009
25	monitor	10000ms	ok (rc=0)	complete	Tue Oct 13 16:31	190ms	0ms	Tue Oct 13 16:31:33 2009

The status bar at the bottom indicates "Connected to 127.0.0.1 (Simple Mode)".

# Sandboxes with graphics!



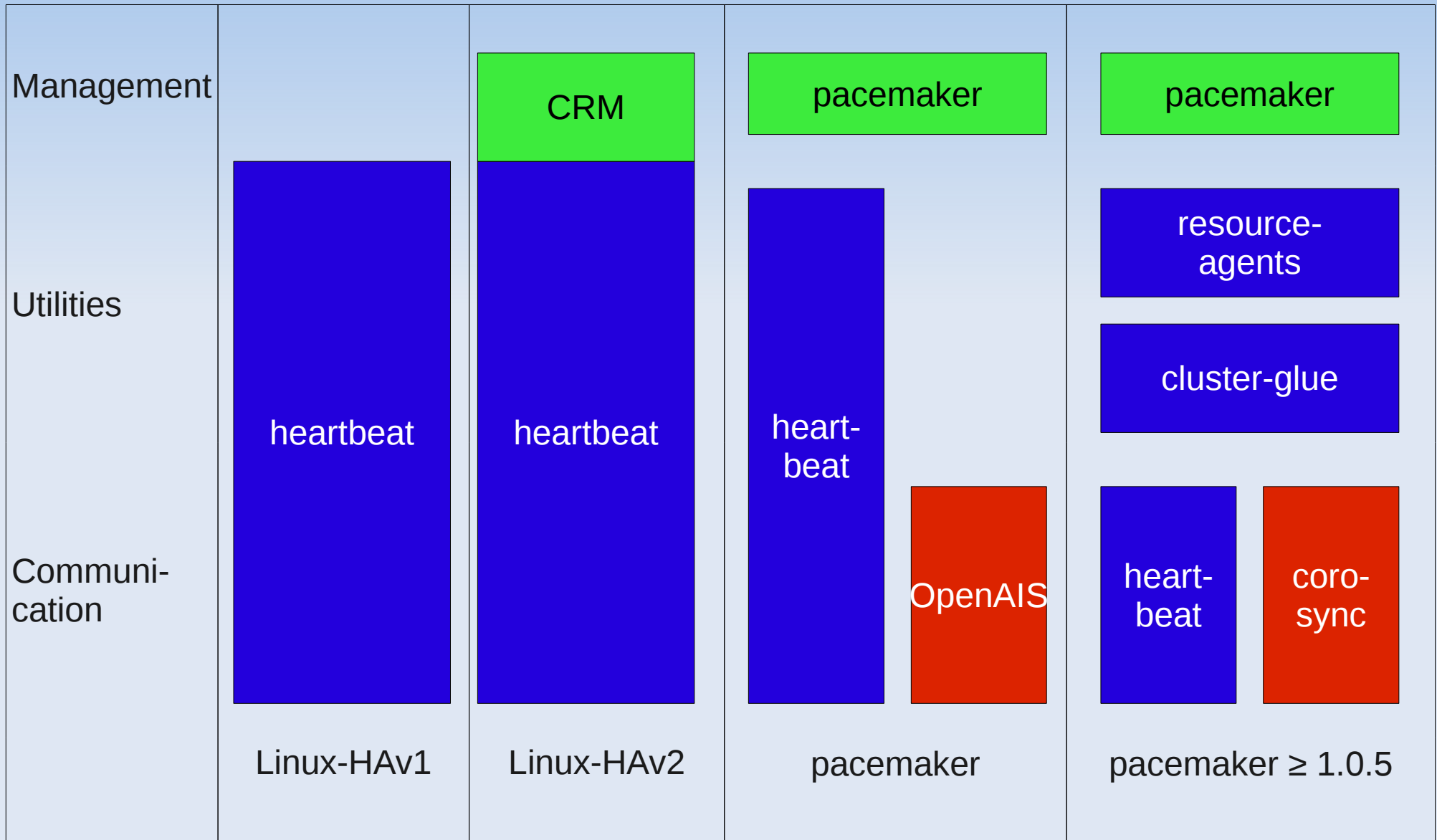
# The development goes on

- OpenAIS offers much more as needed by pacemaker. The communication stack is forked into a separate project called corosync.
- corosync is responsible for the communication in the cluster.
- OpenAIS takes care for all "higher" services in a cluster.
- Simple clusters with pacemaker only need corosync.

# Decomposition of heartbeat

- With pacemaker 1.0.5, heartbeat is decomposed into three projects:
  - `cluster-glue` includes all necessary utilities.
  - The package `resource-agents` combines all resource agents, the interface between pacemaker and the application's binaries.
  - In `heartbeat` (version  $\geq 3.0.2$ ) remains, what is left from the project.

# HA-cluster with Linux





# Master Control Process

- `pacemaker` in version 1.1 got a Master Control Process. So it can be started outside from `corosync` though the `init`-system.
- Of course, `corosync` has to be running.
- Causes less trouble.
- Is more stable.

# DRBD Managment Console

- From company Linbit.
- Written in Java.
- Originally ment to manage DRBDs
- Now full blown cluster management tool:
  - Nodes, resources, constraints
  - Graphical representation of the configuration
  - Emphasis on virtual machines: Config & console.

# Utilizations in pacemaker 1.1

- Every node has utilizations like CPUs, RAM configured.
- Every resource has utilizations configured.
- The cluster manager takes care that resources only run on nodes with enough utilizations.
- TODO: Make resource utilizations dynamic.

# Role Based Access Model

- From version 1.1 administrators access to resources can be limited:
  - You can only destroy your own resources.
- Complete RBAC model implemented.
- Think about a provider offering HA services.

# External Quorum Providers

- Quorum descisions are quite simple up to now.
- `corosync` can use external quorum providers.
- `cman` is a possible external provider.
  - For now it is the only one.

# Beering

- Lavandevil  
Schusterusstrasse 3

U Richard Wagner Platz: U7 / M45

Thank you very much for your attention!

Questions?