

Integrated monitoring of the ATLAS online computing farm



S. Ballestrero · F. Brasolin · D. Fazio · C. Gament · C.J. Lee · D.A. Scannicchio · M.S. Twomey

Why

The online farm of the ATLAS experiment at the LHC consists of nearly 4000 PCs with various characteristics and the status and health of every host must be constantly monitored to ensure the correct and reliable operation of the whole online system.

The monitoring is the first line of defense: it should not only promptly provide alerts in case of failure but, whenever possible, warn of impending issues.

How

Two different monitoring systems have been combined to get a complete picture of the system:

- ★ Icinga 2 for active checks and alerting
 - + it replaced Icinga (end 2015) which in turn replaced Nagios v3 (2014/2015)
 - + it improved scalability and has native support for load balancing
 - ▶ it currently handles ~ 4000 hosts and ~ 83200 checks performed with different time intervals ranging between 5 minutes and 24 hours according to the need
 - + two servers are configured with the built-in High Availability Cluster feature
- ★ Ganglia for performance data useful for debugging
 - + historical data are stored
 - + high scalability and good data visualization
 - + Ganglia Monitoring Agent (gmond)
 - ▶ runs on each node and gathers data every 20 seconds
 - ▶ some parameters useful for alerting too
 - ▶ helps reducing the active Icinga 2 checks

In addition Ganglia and Icinga 2 have been integrated with other data sources, such as SNMP for system information and IPMI for hardware health.

What

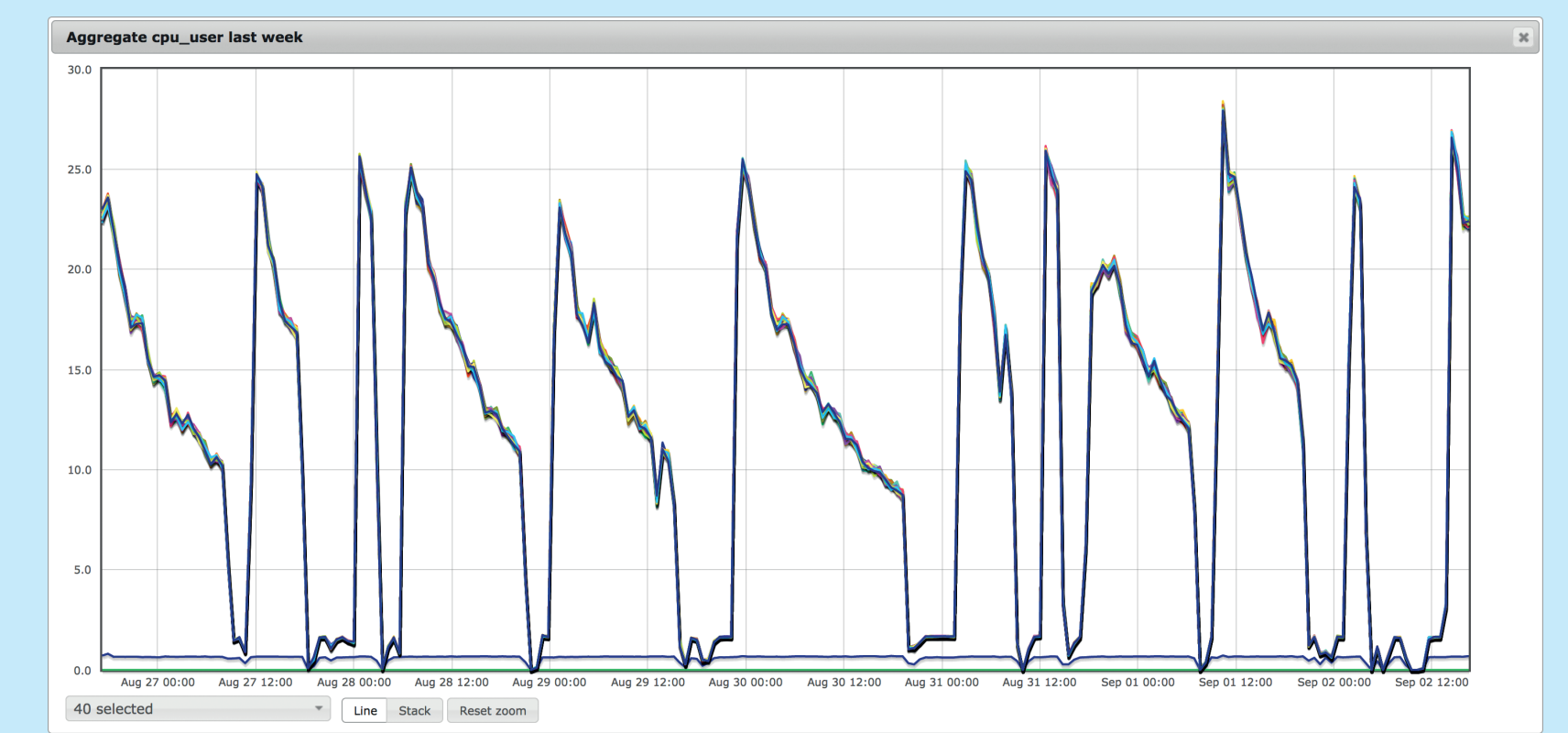
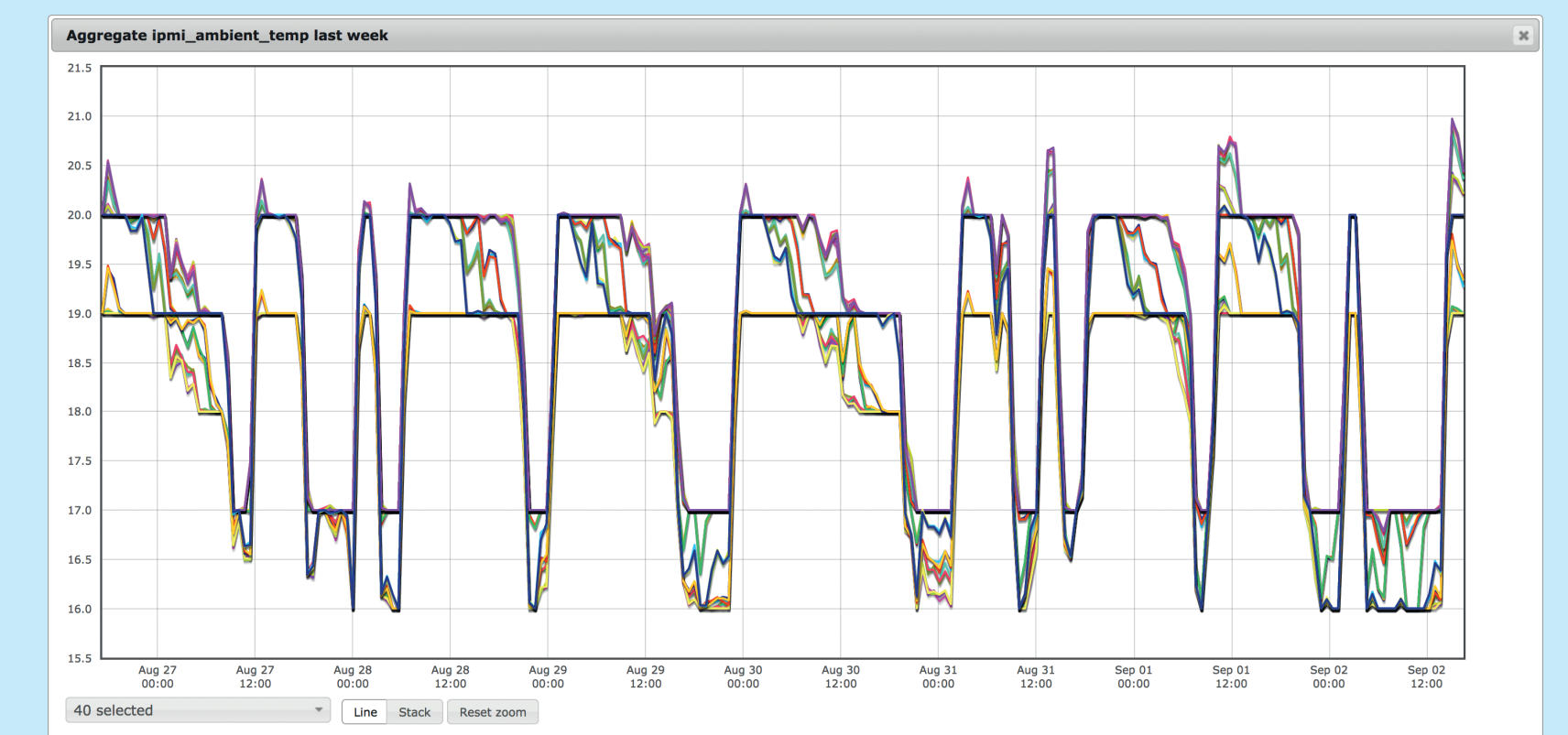
The monitoring system should be able to check up to 100000 health parameters and provide alerts on a selected subset.

For each node various (~20) hardware and system parameters are monitored:

- ★ hardware:
 - + disk raid status, temperature, fan speed, power suppliers, current and voltages, firmware version...
- ★ system:
 - + cpu, memory and disk usage
 - + configuration aspects as kernel version, pending updates
 - + network: interfaces status, speed and bonding (if present)
 - + services as ntp, http, ssh, mailq

The huge number of nodes and variety of configuration and settings required to set up an automatic way for creating the Icinga 2 configuration files:

- ★ general Icinga 2 templates have been defined
- ★ the in-house Configuration Database (ConfDB) stores the information related to hardware, OS, server types (netbooted or localboot) and specific settings for the Icinga 2 checks, e.g. overriding default thresholds
- ★ SQL queries and regexp, exploiting the node naming convention, are used to extract data and select the right template for the host



This guarantees the full coverage of all the nodes, a per-type configuration and per-host tuning.

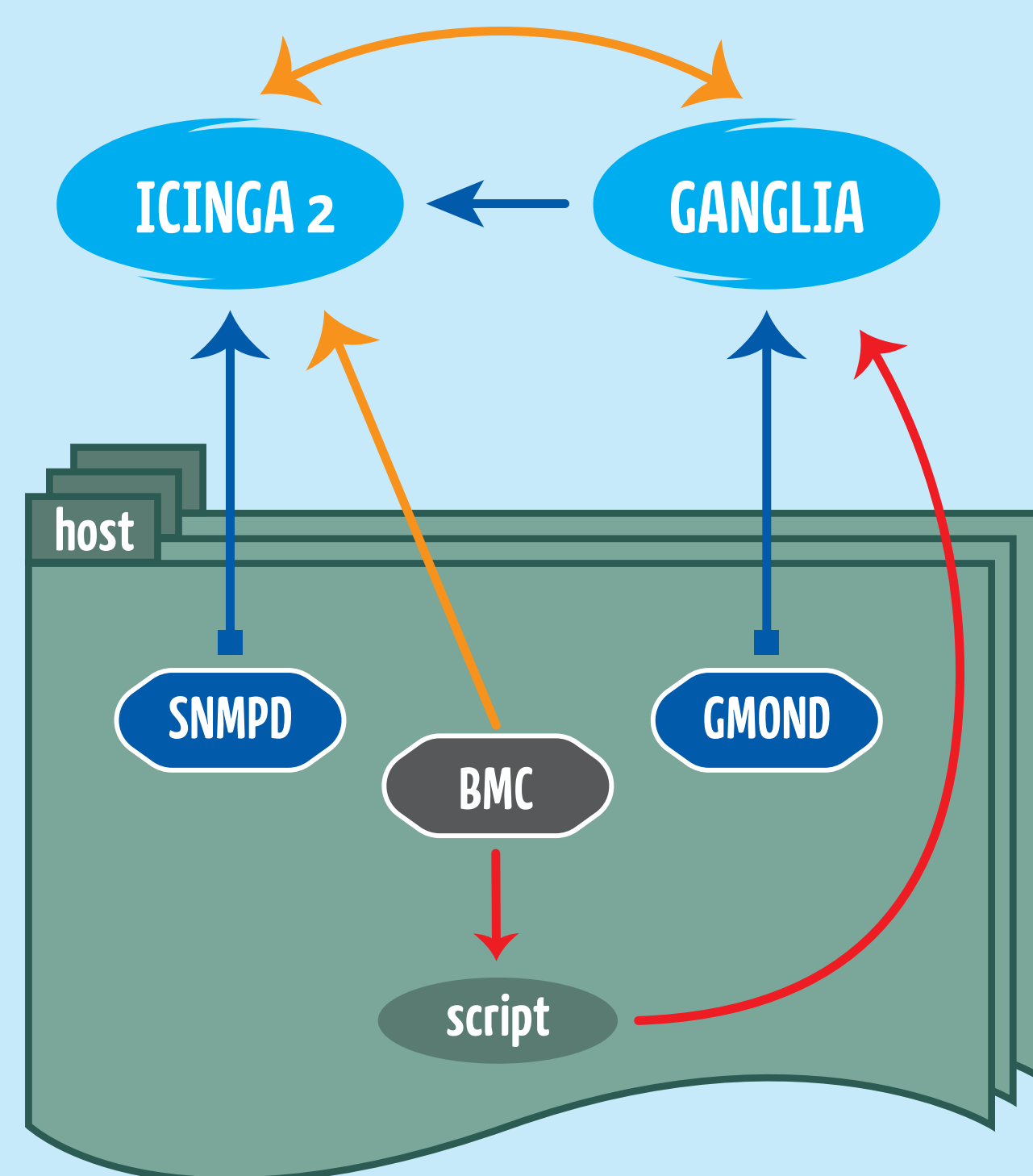
A big effort has been put into customising the monitoring of all the significant IPMI sensor information provided by the Baseboard Management Controllers (BMC) of the PCs.

The two Ganglia plots show respectively the ambient temperature and CPU usage of each of the 40 nodes in a rack:

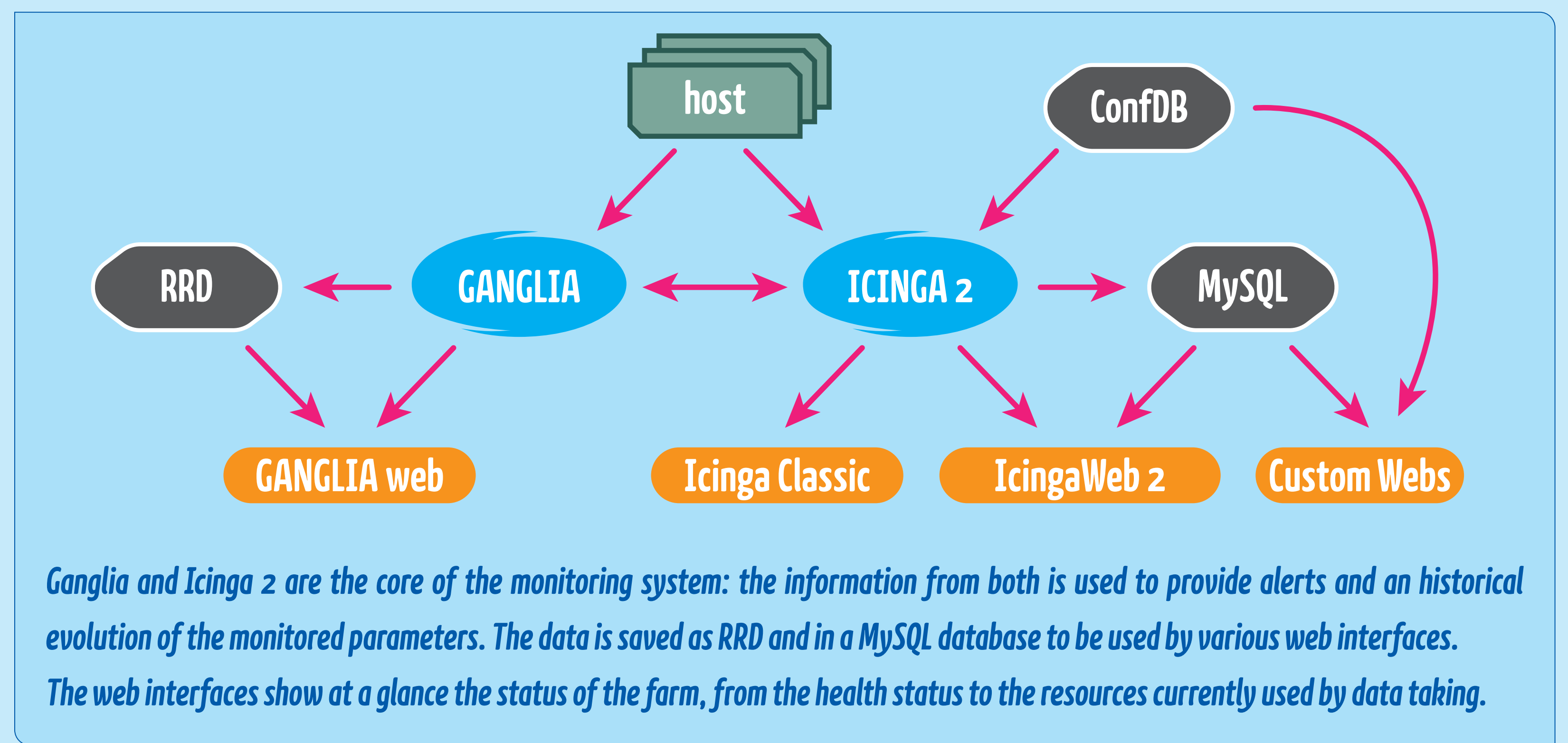
- ▶ the temperature is sent to Ganglia by the IPMI sensors on the nodes
- ▶ the CPU usage is sent to Ganglia by the Ganglia Monitoring Agent running on each node

Monitoring system

- DATA FLOW:**
- SNMPD**
data polled by Icinga 2 (blue lines)
 - GMOND**
data pushed to Ganglia, Icinga 2 can poll Ganglia (blue lines)
 - BMC**
Icinga 2 polls critical data from BMC, stores copy to Ganglia (orange lines)
 - script polls the host's BMC, pushes complete data to Ganglia (red lines) (no info if host is down)



The arrows indicate the data flow



Ganglia and Icinga 2 are the core of the monitoring system: the information from both is used to provide alerts and an historical evolution of the monitored parameters. The data is saved as RRD and in a MySQL database to be used by various web interfaces. The web interfaces show at a glance the status of the farm, from the health status to the resources currently used by data taking.

Farm overview

GROUPS	TOTAL	ONLINE	OFFLINE	MAINT	RESERVED	Comments
@ Gateways	6	6	0	0	0	Full emergency info
@ Hosts	6	6	0	0	0	Full emergency info
@ CoreServers	36	36	0	0	0	Dedicated to Swift only
@ SCRs	30	27	2	1	0	Disabled
@ LFS	54	54	0	0	0	Full emergency info
@ ONL	2	2	0	0	0	Full emergency info
@ ONL	34	34	0	0	0	Full emergency info
@ MON	40	40	0	0	0	Full emergency info
@ CAL	32	32	0	0	0	Full emergency info
@ RDS	102	102	0	0	0	Full emergency info
@ SFO	2	2	0	0	0	Full emergency info
@ HLTSTV	2	2	0	0	0	Full emergency info
@ TRU	2892	2817	39	36	0	Full emergency info
@ CTP	5	5	0	0	0	Full emergency info
@ BST	5	5	0	0	0	Full emergency info
@ RMON SRVs	2	2	0	0	0	Full emergency info
@ NETWORK	81	81	0	0	0	Full emergency info
@ VAL	4	4	0	0	0	Full emergency info
@ SEC	4	4	0	0	0	Full emergency info
@ PUB	13	13	0	0	0	Full emergency info
@ OCS	186	186	0	0	0	Full emergency info
@ MU-CALSRV	290	285	0	0	0	Full emergency info
@ DETECTOR	258	241	0	0	17	Full emergency info
@ SWITCH	6	6	0	0	0	Full emergency info
@ SIMPL	18	18	0	0	0	Full emergency info
@ TestSquare	66	66	0	0	0	Full emergency info
@ OTHERS	4115	3969	75	40	23	Full emergency info

A custom web page shows an overview of the status of the whole farm.

A "cross-consistency check" custom web page provides information from Icinga, ConfDB and the ATLAS Run status: it displays a quick overview of the resources not in use by ATLAS because of some ongoing hardware or software intervention.

hostname	Icinga	ConfDB HS	TDAQ Part	Issues	
pc-tdq-tpu-03010	slow	ack=1	203	maint,sysadmin	oks/0 out (set Icinga maint)(ok to intervene)
pc-tdq-tpu-03022	up	ack=0	2011	prod,tdaq	oks/1 out
pc-tdq-tpu-04019	up	ack=0	2011	prod,tdaq	oks/1 out

Conclusions

The new monitoring system based on Icinga 2 and Ganglia provides the required information and alert notifications. Our experience shows that a comprehensive and robust monitoring system allows to prevent a lot of hardware and software issues in advance, before they become critical for the ATLAS data taking and the security of the whole system.