





CERN Cloud Infrastructure update

Daniel Abad on behalf of the Cloud Infrastructure Team
HEPiX Autumn 2019 Workshop

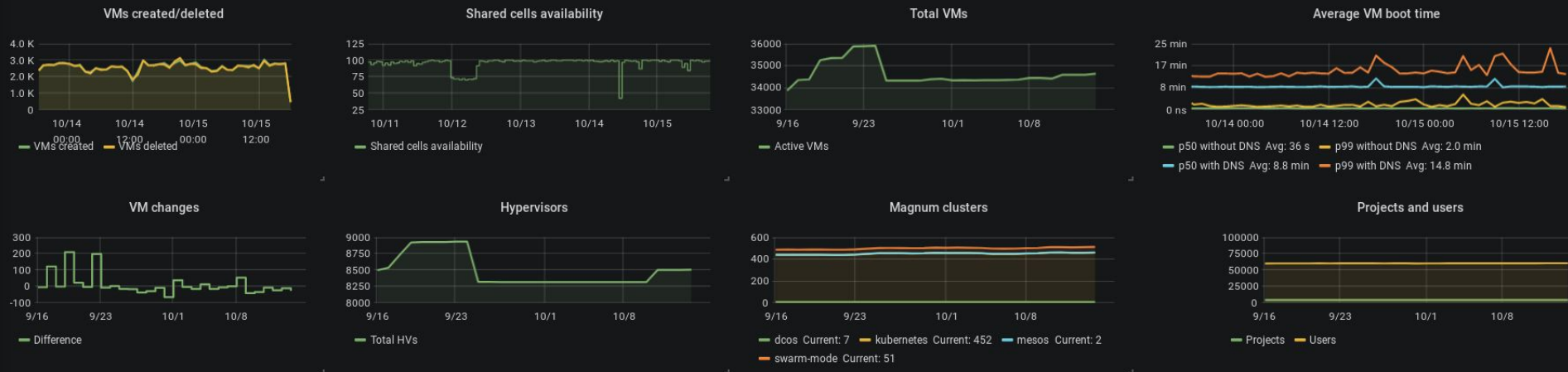
Cloud resources



Openstack services stats



Resource overview by time



Development areas going forward

- Spot Market / Pre-emptible instances
- Software Defined Networking
 - Introducing LBaaS this month
- Magnum rolling upgrades
- Enrolling all 15K servers in Ironic
 - Containers on Bare Metal
- Evaluation of OpenStack Watcher for resource optimization
- Upgrading to Stein release

<https://techblog.web.cern.ch/techblog/>

OpenStack Days CERN, 27MAY2019





<https://indico.cern.ch/event/776411/>

L1TF (Reboot campaign again)

- Disabling SMT OFF in the BIOS makes the difference respect to acting at kernel level. We can still:
 - Overcommit resources and provide the same performance as before
 - Have servers patched for L1TF
 - D. Giordano on “[CPU steal and L1TF mitigations](#)”
- Plan to fully patch L1TF
 - Continue the reconfiguration of the remaining nodes that still have SMT ON
 - Set SMT OFF in the BIOS
 - A reboot campaign of the affected HVs was announced, VMs’ users contacted and finally performed in May

Container Orchestration - OpenStack Magnum

An OpenStack API service that allows creation of container clusters

- Use your OpenStack Keystone credentials
 - You choose your cluster type
 - Multi-tenancy
 - Quickly create new clusters with advanced features
-
- Deprecating Mesos/DCOS  MESOS  DC/OS
 - Today, the CERN Cloud hosts >500 Kubernetes clusters



MAGNUM
an OpenStack Community Project



Bare metal Provisioning - OpenStack Ironic



IRONIC

an OpenStack Community Project

Integrated OpenStack program which aims to provision bare metal machines instead of virtual machines

- In production at CERN since 2018
- Consolidate accounting & bookkeeping
- Completes our service offering
- All new hardware is enrolled using Ironic
- +3700 nodes managed by Ironic
- Testing enrolment of existing hardware

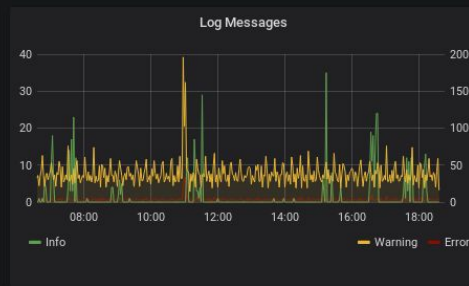
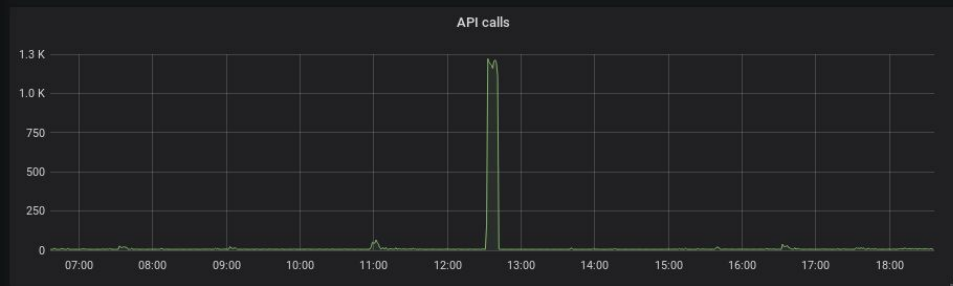
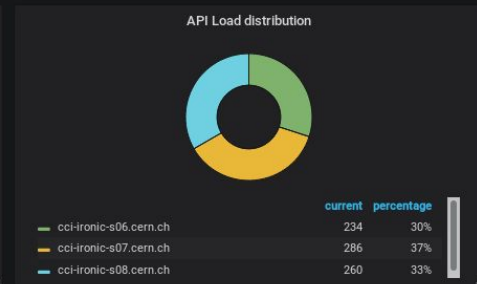
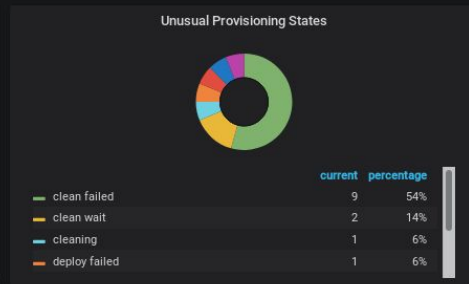
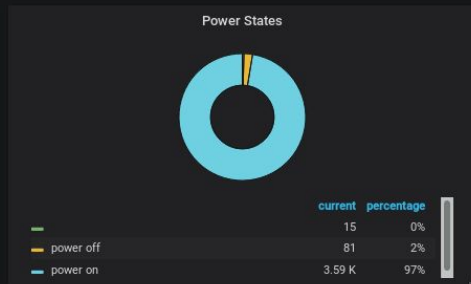
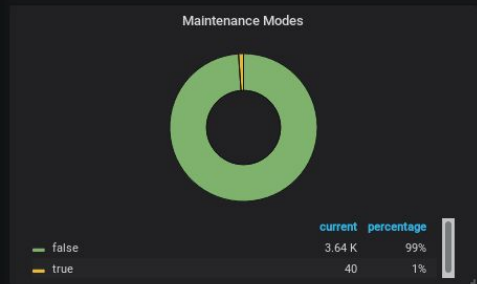
Why Bare-Metal Provisioning? (1)

- VMs not suitable for 100% of our use cases
 - Benchmarking, storage nodes, boot strapping, critical network equipment, specialised network setups, HPC clusters, ...
- Complete our service offerings
 - Physical nodes (in addition to VMs and containers)
 - OpenStack UI as the single pane of glass
- Simplify hardware provisioning workflows
 - For users: `openstack server create/delete`
 - For procurement: initial on-boarding, server re-assignments



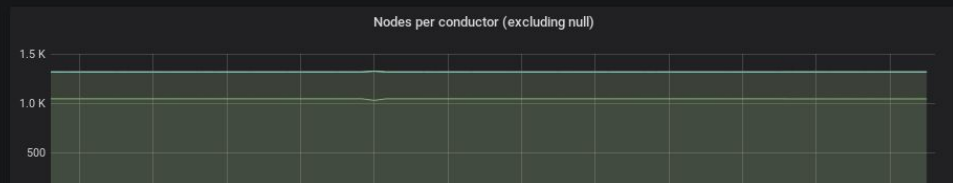
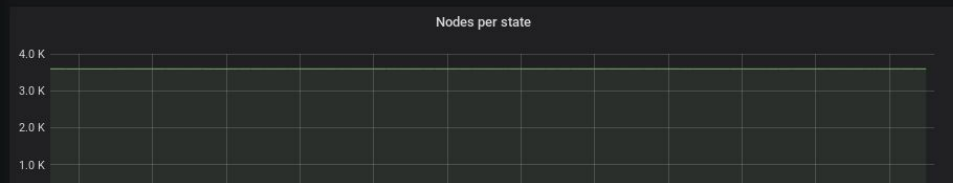
Integrating Ironic into the CERN's Private Cloud Service, HEPiX, Madison 2018

3



Rally Tests

Time	bm_test	bm_qa	bm_001
6:00 PM	0	0	-
5:00 PM	0	0	1
4:00 PM	0	0	1
3:00 PM	0	0	0
2:00 PM	0	0	0
1:00 PM	0	0	0



Getting physical resources in 2011

Public Procurement Purchase Model

Step	Time (Days)	Elapsed (Days)
User expresses requirement		0
Market Survey prepared	15	15
Market Survey for possible vendors	30	45
Specifications prepared	15	60
Vendor responses	30	90
Test systems evaluated	30	120
Offers adjudicated	10	130
Finance committee	30	160
Hardware delivered	90	250
Burn in and acceptance	30 days typical 380 worst case	280
Total		280+ Days

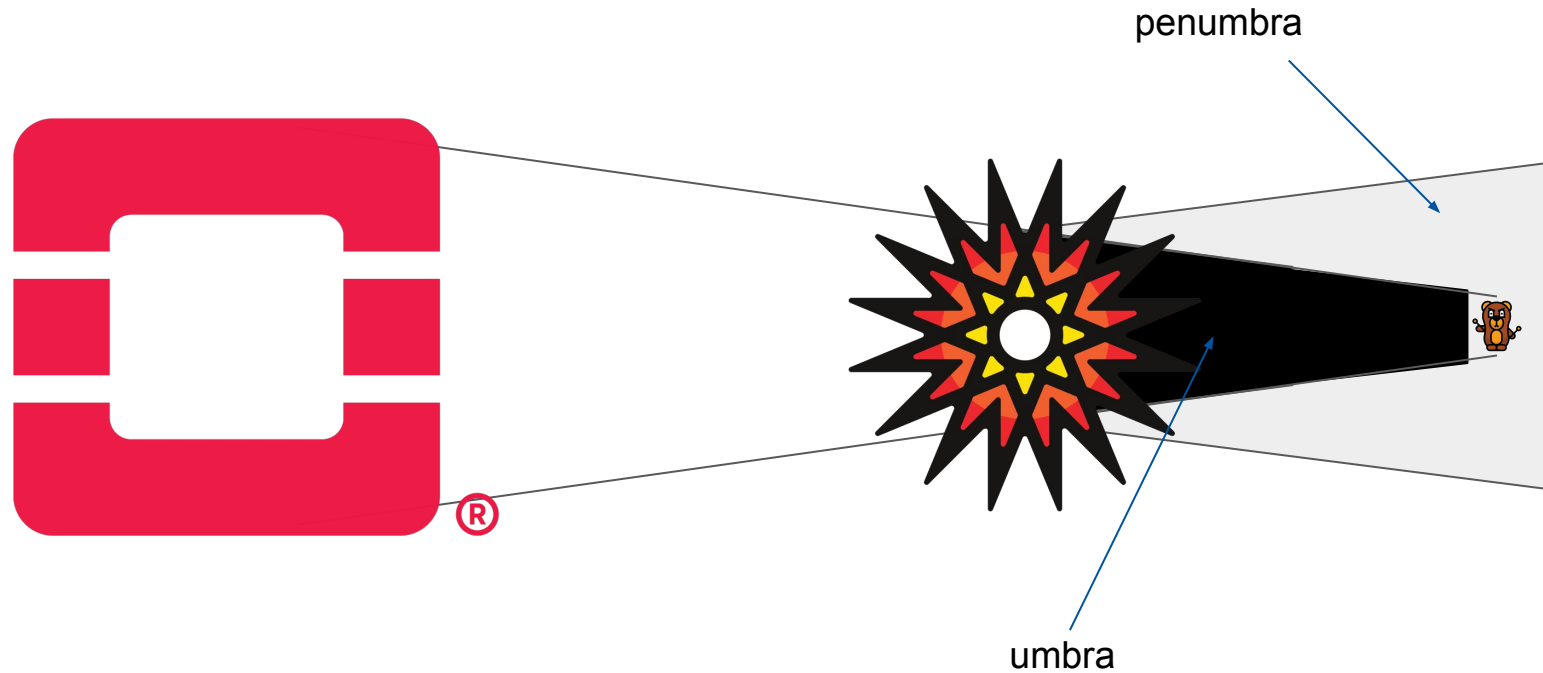
Getting physical resources in 2019

```
dabadaba@aiadm $ openstack server create \  
  --image "CC7 - x86_64 [2019-09-30]" \  
  --flavor p1.cd6839123.S513-H-IP71 \  
  my_physical_instance
```

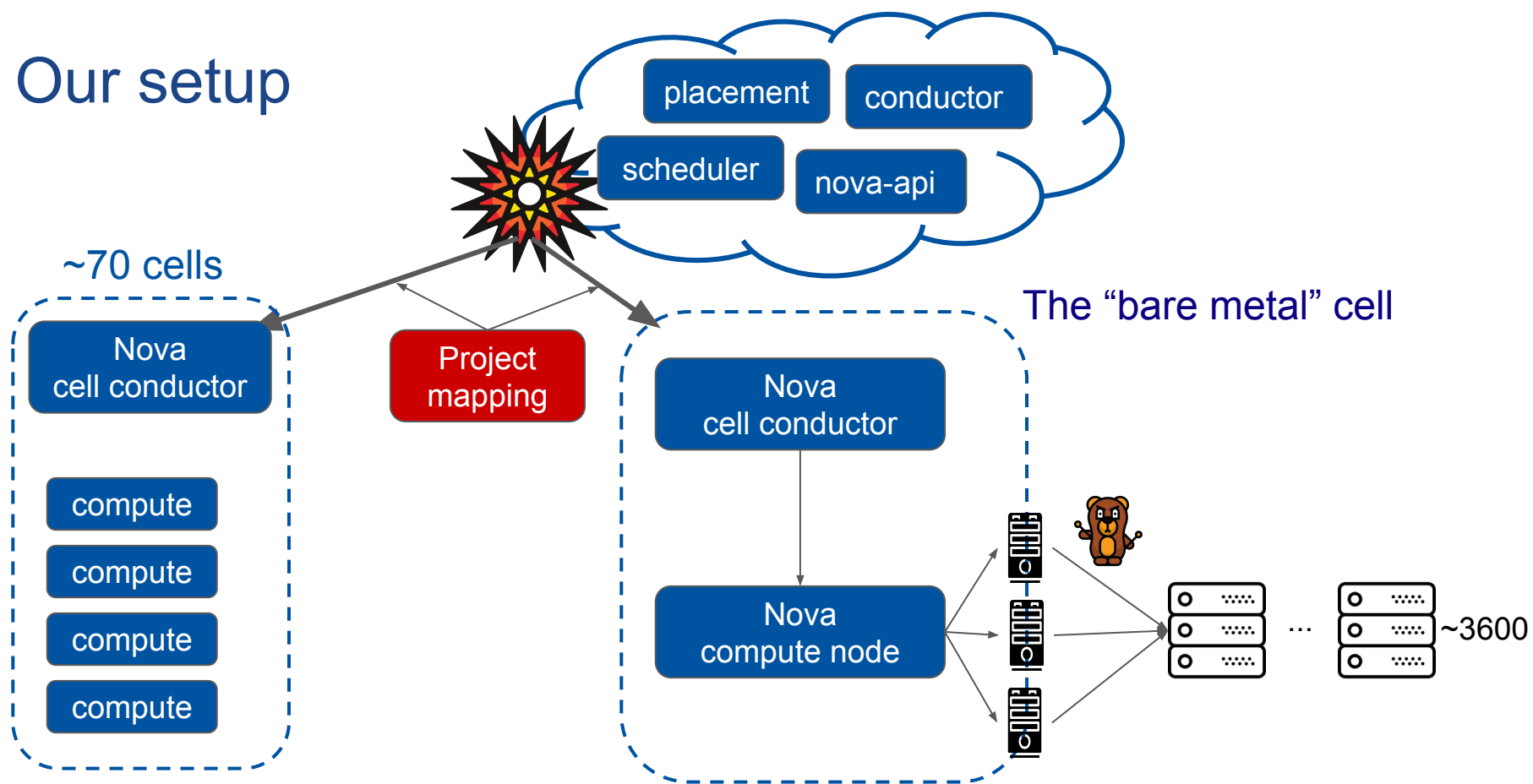
15 minutes later

```
dabadaba@aiadm $ ssh root@my_physical_instance  
root@my_physical_instance #
```

Supernova eclipse



Our setup

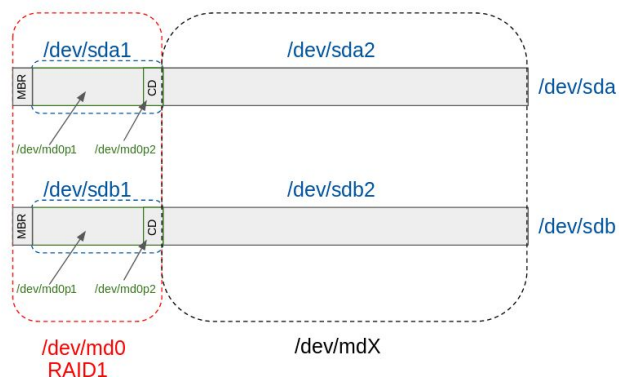


Scalability issues

- Started with ~1500 nodes
 - (And only one level 2 compute node)
 - Overcome and now working for >3600 nodes
- Mainly synchronization between nova and Ironic
 - Upgrade to placement on nova
 - Placement designed for virtual resources (VMs, GPUs, etc.)
 - Tweaked some params for the bare metal cell:
- <https://techblog.web.cern.ch/techblog/post/nova-ironic-at-scale/>

Software RAID

- Sep 2018 - Proposal: [Downstream Ironic Software RAID support at CERN](#)
- Dec 2018 - Specification:
 - <https://specs.openstack.org/openstack/ironic-specs/specs/not-implemented/software-raid.html>
 - <https://review.opendev.org/#/c/624413/>
- Jan-Jun 2019 - Implementation
 - https://techblog.web.cern.ch/techblog/post/ironic_software_raid/



```
ironic node-set-target-raid-config  
my-baremetal-node-1379  
{  
  "logical_disks": [{  
    "raid_level":  
    "1",  
    "size_gb": "MAX",  
    "is_root_volume": true  
  }]}'
```

Operational issues

- Adoption of existing hardware
 - ~8000 hypervisors
 - Testing the strategy, high risk
- Mix of Physical resources:
 - Nodes outside of OpenStack
 - Hypervisors
 - End-user nodes (Storage, Windows, etc.)
- Different procedures for interventions
 - Coordination
 - Human communication

Inventory of Hardware resources: a CMDB?

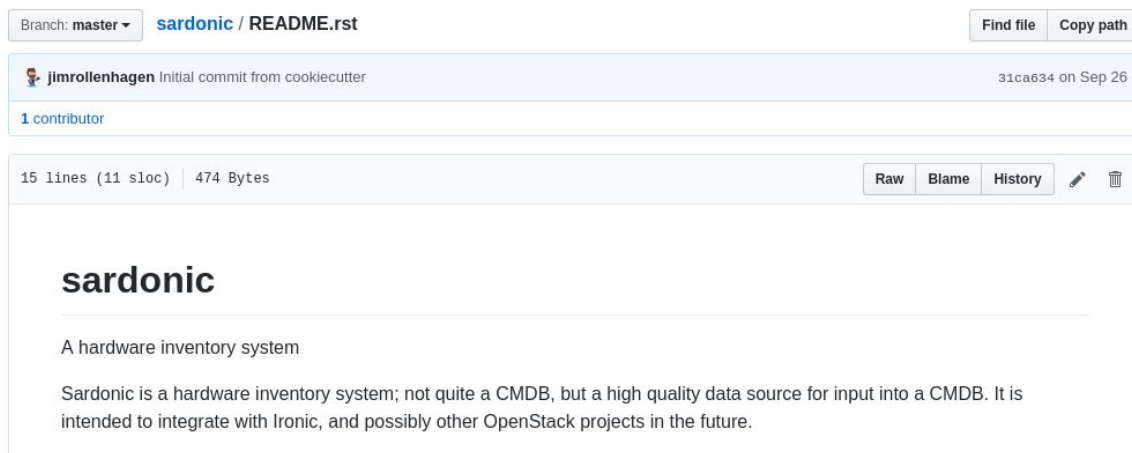
The image shows a large, multi-column spreadsheet or data table. The table is filled with text and numbers, representing an inventory of hardware resources. Several rows and columns are highlighted in green, yellow, and blue, indicating different categories or statuses of the resources. A prominent red horizontal line is drawn across the bottom portion of the table, possibly indicating a summary or a specific section. The overall appearance is that of a detailed data management tool, such as a CMDB (Configuration Management Database).

CMDB: main use cases

- Single source of truth
- Inventory of servers and components
 - Accounting
 - Including temporary “leases”
- Upgrade campaigns
 - Hardware: SSDs, CPUs, etc.
 - Software: kernels, etc.
- Warranty information
 - inventory of servers and components

Sardonic - an OpenStack Proposal

Proposal: [Ideas and Requirements for a "Hardware Inventory" Functionality in OpenStack Ironic](#)



The screenshot shows a GitHub repository page for the 'sardonic' project. At the top, it indicates the current branch is 'master' and the file being viewed is 'sardonic / README.rst'. There are buttons for 'Find file' and 'Copy path'. Below this, a commit by 'jimrollenhagen' is shown, with the message 'Initial commit from cookiecutter' and the commit hash '31ca634' dated 'on Sep 26'. It lists '1 contributor'. The file statistics show '15 lines (11 sloc) | 474 Bytes'. There are buttons for 'Raw', 'Blame', and 'History', along with edit and delete icons. The main content of the README is as follows:

sardonic

A hardware inventory system

Sardonic is a hardware inventory system; not quite a CMDB, but a high quality data source for input into a CMDB. It is intended to integrate with Ironic, and possibly other OpenStack projects in the future.

<https://github.com/openstack/sardonic>

CMDB: solutions evaluated

Solution	Open Source	License	Platforms	Business model	Language	Last release
I-doit	Yes, no repo	not specified	Linux/Windows	Free + Pro (differences)	PHP5 + JS	1.12.2 (2019-04-01)
CMDbuild	Yes	AGPL v3.0	Linux/Windows	Free + Pro (differences)	Java + Ajax	3.0.0 (2019-04-12)
SnipeIT	Yes	AGPL v3.0	All	Free + Pro (differences)	PHP + JS	4.6.15/5.0.0-beta-2 (2019-03-20)
SysAid	No	/	Windows (All)	paid	???	18.1 (2019-11-04)
Spiceworks	No	/	Windows only	freeware	Ruby (2006)	7.5.00107 (2018-03-21)
GLPI	Yes	GPL v2	All (linux mainly)	libre	PHP + JS	9.4.2 (2019-04-11)
NetBox	Yes	Apache 2.0	Linux	libre	Python 3.5	2.5.10 (2019-04-08)
xCAT	Yes	EPL 1.0	Linux	libre	Perl	2.14.6 (2019-03-29)
openDCIM	Yes	GPL v3	Linux	libre	PHP + JS	19.01 (2019-03-04)
Ralph	Yes	Apache 2.0	Linux	libre	Python 3.6 + JS	20190410.2 (2019-04-10)

CMDB: solutions evaluated

Solution	Supported DBs	Schema flexibility	API	GUI
i-doit	MySQL/MariaDB	???	JSON-RPC	Web
CMDBuild	PostgreSQL only	Yes	REST, SOAP	Web (Framework)
SnipeIT	MySQL/MariaDB	At least partial	JSON REST	Web
GLPI	PostgreSQL/MySQL(InnoDB)	???	REST	Web
NetBox	PostgreSQL	???	REST	Web
xCAT	SQLite, MySQL/MariaDB, PostgreSQL	???	REST	Web
openDCIM	MySQL	???	REST	Web
Ralph	MySQL	???	REST (Django)	Web

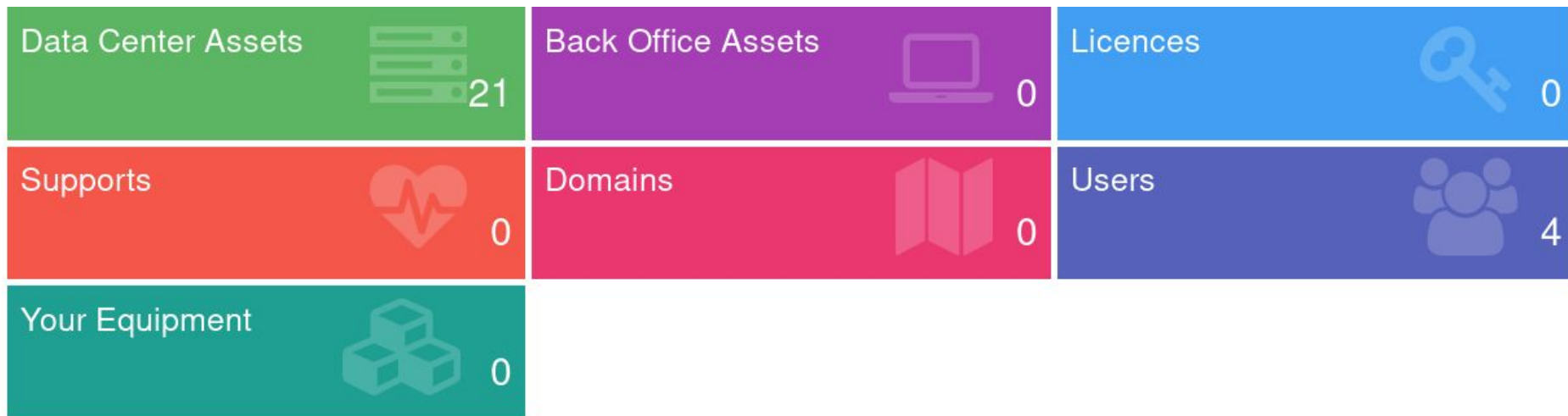
CMDB: solutions evaluated

Solution	Supported Linux distributions	Installation from repo/ Automatic install
i-doit	Debian 8/9, Ubuntu 16.06/18.04, RedHat (RHEL) 7, SUSE (SLES) 12 SP2, SP3, SP4	Yes, with bash script <code>idoit-install</code>
CMDBuild	Any	No
SnipeIT	Any that supports a LAMP stack or equivalent	Partial, web installation
GLPI	Any that supports a LAMP stack or equivalent	Yes, web or cli
NetBox	Ubuntu 18.04, CentOS 7.5	No
xCAT	Ubuntu Server LTS 18.0.4, RHEL 7.6 , SLES 12.3	Yes, with the bash script <code>go-xcat</code> or configuring repositories
OpenDCIM	Any that supports a LAMP stack or equivalent	Partial, web installation
Ralph	Ubuntu 14.04 / 18.04	Install from repo; manual configuration of web stack

CMDB: solutions tested

- NetBox
 - Complex installation
 - Mainly IP address management (IPAM) driven
- Ralph
 - Written in Python/Django
 - Better support for custom fields
 - Easy to configure
- xCAT
 - Does not have a GUI

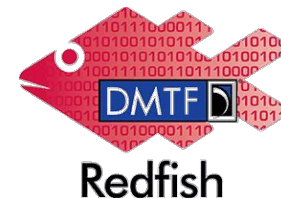
Ralph: Testing with real data



Ralph: Testing with real data

<input type="checkbox"/>	Hostname	Status	Barcode	Model	SN	Invoice date	Invoice number	Location	Service env	Configuration path	Security scan	SCM status
<input type="checkbox"/>	(None)	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	a legit serial number v4	(None)		1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]51.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]59-4	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]37.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]59-3	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]73.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]59-2	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]77.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]59-1	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]83-3	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]83-2	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]78-3	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]78-2	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]78-1	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED].cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 2	[REDACTED]69-3	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]642-b.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]642-B	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]642-a.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]642-A	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]641-d.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]641-D	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]641-c.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]641-C	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]641-b.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]641-B	(None)	(None)	1	service1 - environment1	(None)	No scan	-
<input type="checkbox"/>	[REDACTED]641-a.cern.ch	new	(None)	[Data Center] Manufacturer1 [REDACTED] 5	[REDACTED]641-A	(None)	(None)	1	service1 - environment1	(None)	No scan	-

Redfish integration

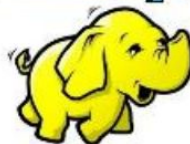


- TL;DR IPMI's “successor”
- Driven by DMTF Scalable Platforms Management Forum
 - **Promoters:** Broadcom, Dell, Emerson, Hewlett-Packard, Intel, Lenovo, Microsoft, Supermicro, VMWare
 - **Supporters:** AMI, Oracle, Fujitsu, Huawei, Mellanox, Seagate
- An open industry standard specification that provides simple, modern and secure management of scalable platform hardware
- RESTful interface over HTTPs in JSON format
- A secure, multi-node capable replacement for IPMI-over-LAN

The next 5 years?

- LHC to Run 4
 - Computing must not limit the physics
- Open Infrastructure
 - OpenStack is a key part but lots of others too
 - Ceph, Tungsten Fabric, Grafana, Puppet, CI/CD, Kubernetes, ...
- Open Source collaboration is the way forward
 - Natural culture fit for sharing and giving back
 - Work with the communities which are open and vibrant

We are not special anymore



Summary

- During the last 10 years, resource management and deployment model changed completely
 - From Virtualization and Server consolidation to Cloud Infrastructure
 - From Bare metal to VMs, to managed Bare metal to Containers
- Continue to adapt the infrastructure to new technologies and requirements
 - Control plane managed by kubernetes
 - New regions
 - SDN
 - Serverless
- Community and industry collaboration has been productive and inspirational for the CERN team

