

When unattended upgrades were really unattended

REX on performance analysis



Introduction

Hello 🖐️ everyone,

I am Nicolas Sterchele, SRE at BackMarket.

(Previously SRE at PeopleDoc, Cloud Engineer at AntVoice...)

@sterchelen on Github/Twitter



Agenda

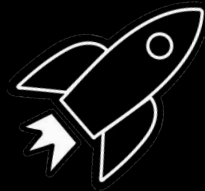
- Current Architecture
- The issue that had an atomic clock precision
- Methodologies
 - This is not my hardware
 - USE method
 - Drill-down analysis
- Wise words and takeaways
- Hey, hey oh... we are hiring :)
- References



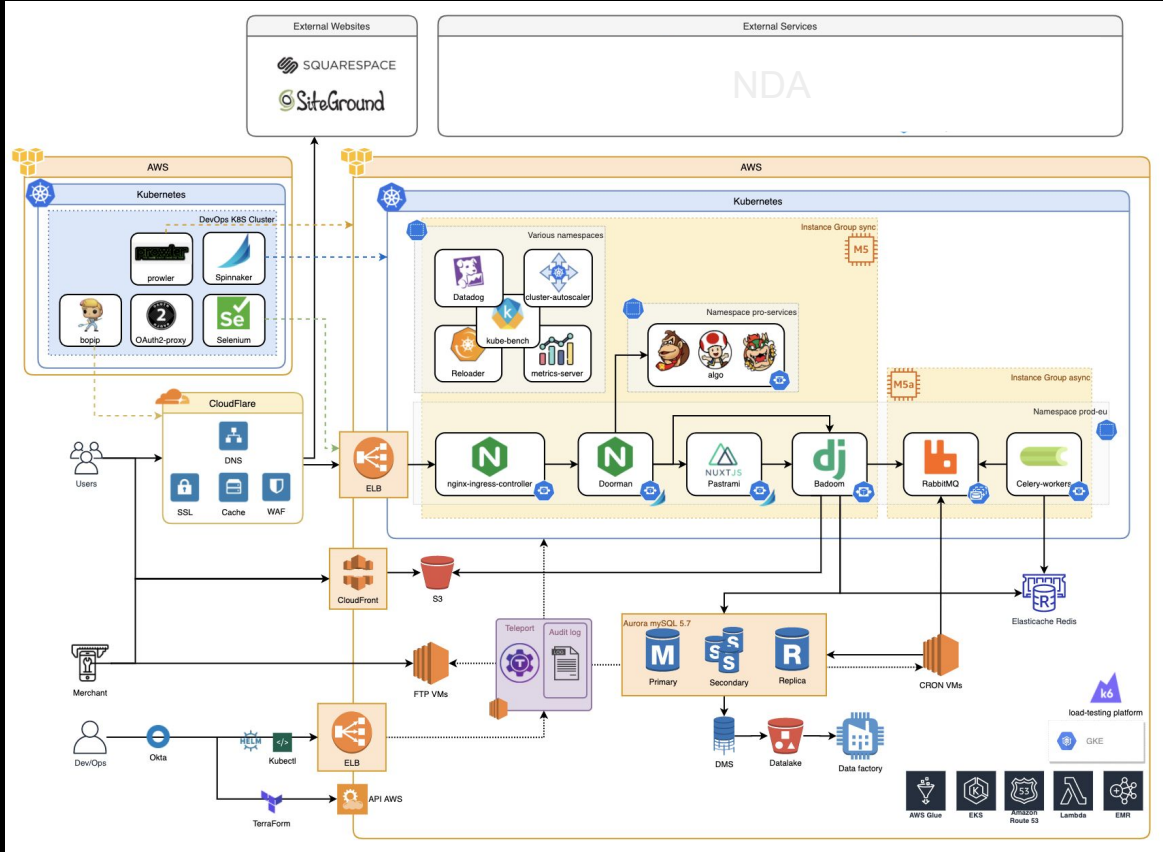
BM Architecture



In a nutshell



BM Architecture



2 AWS regions

Between 40 - 300 K8s nodes

300 - 3000 pods

4 main self-managed clusters

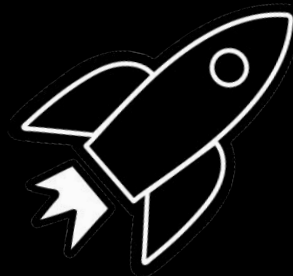
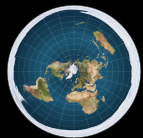
Traffic multiplied by 2 every year

~80 engineers



The issue that had

an atomic clock precision



The issue

Wednesday morning, the 21st of July, the **SRE team** receives a message from an Engineer:

*We had experienced **connectivity issues** to our redis clusters **and** our mysql databases, between 8AM and 9AM.*



The issue

Indeed,

Redis TimeoutError: Timeout reading from socket
SQL Error: (2006, 'MySQL server has gone away')



Since the errors stopped and weren't reaching a critical rate, we decided to monitor and put that on the side.

You know the song about SREs, they always have dozens of things on the fire.

The issue

The very next day, same hour,
same errors, same message
from the Engineer.





Methodologies

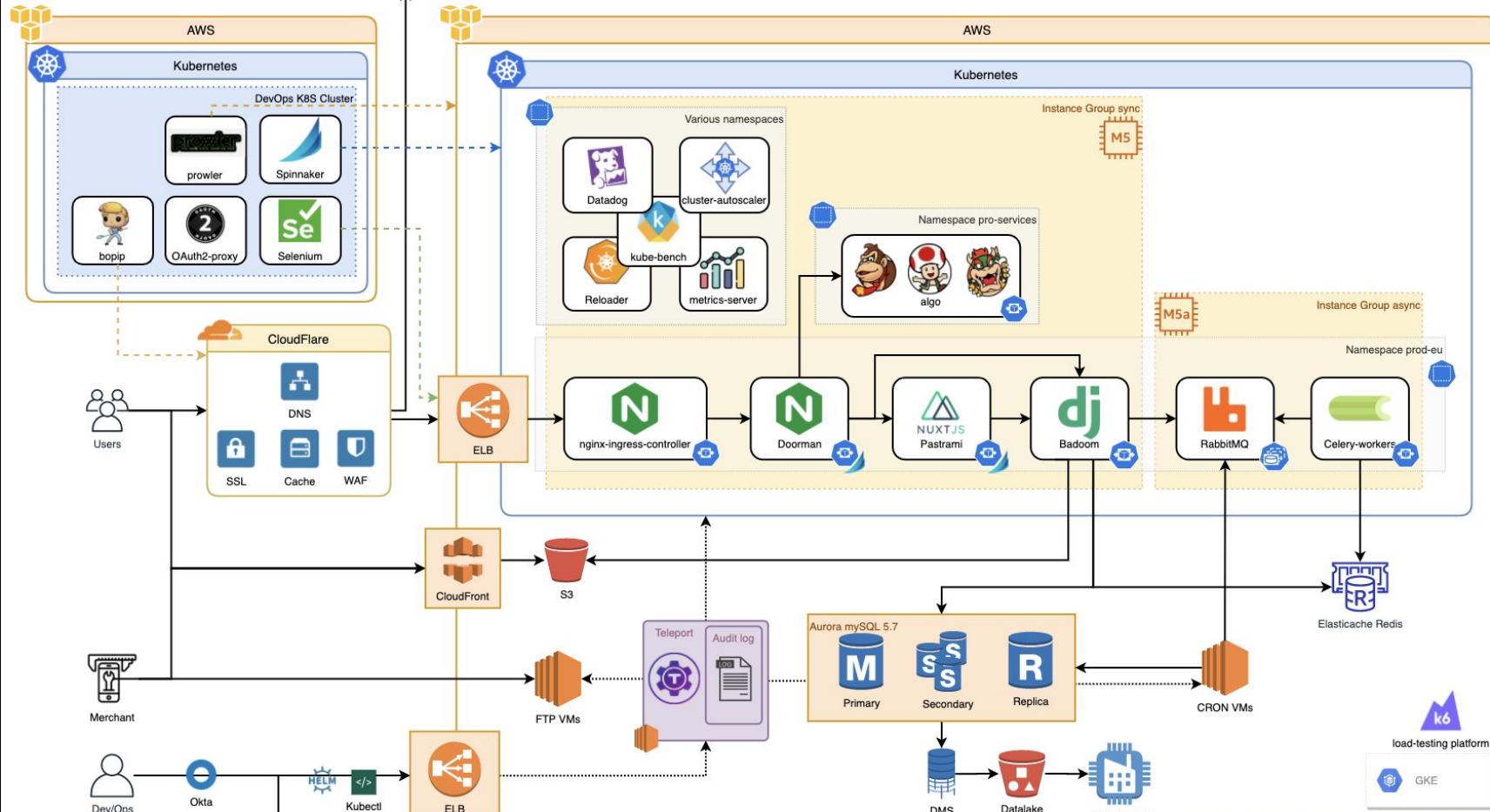
We are scientists, right?



Reminder

SQUARESPACE
SiteGround

NDA



**This is not my
hardware**

Using a public cloud provider is like renting an apartment.
Maintenance is handled by the owner, **AWS** in our case.

We saw on the 21st of July some maintenance events on the
AWS ElastiCache clusters.

Again, we don't have control on what has been applied, that's
why we decided to send an email to the AWS support team.

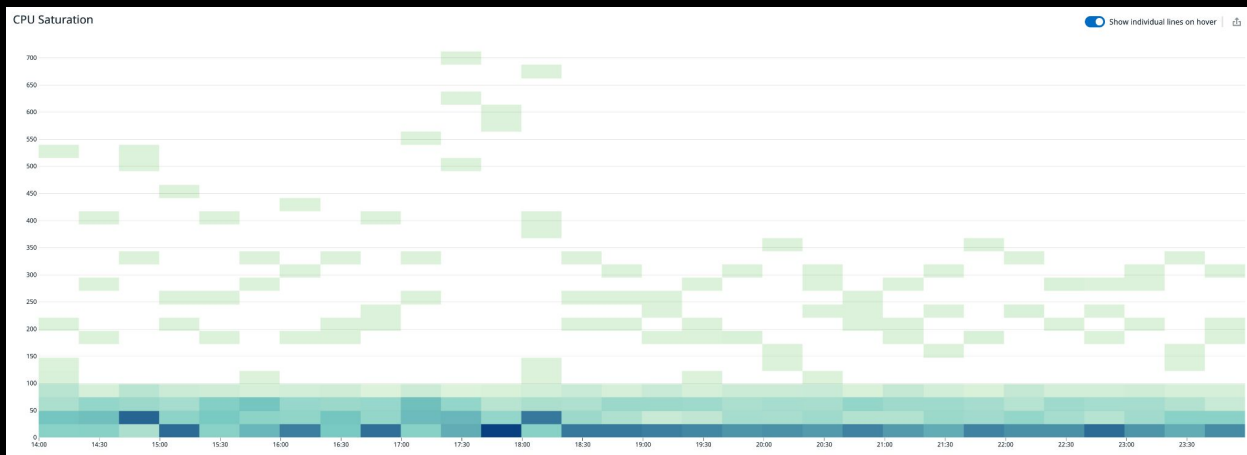
USE Method

The USE method is a methodology that focuses on system resources and can be summarized as [Gregg 01]

For every resource, check utilization, saturation, and errors.

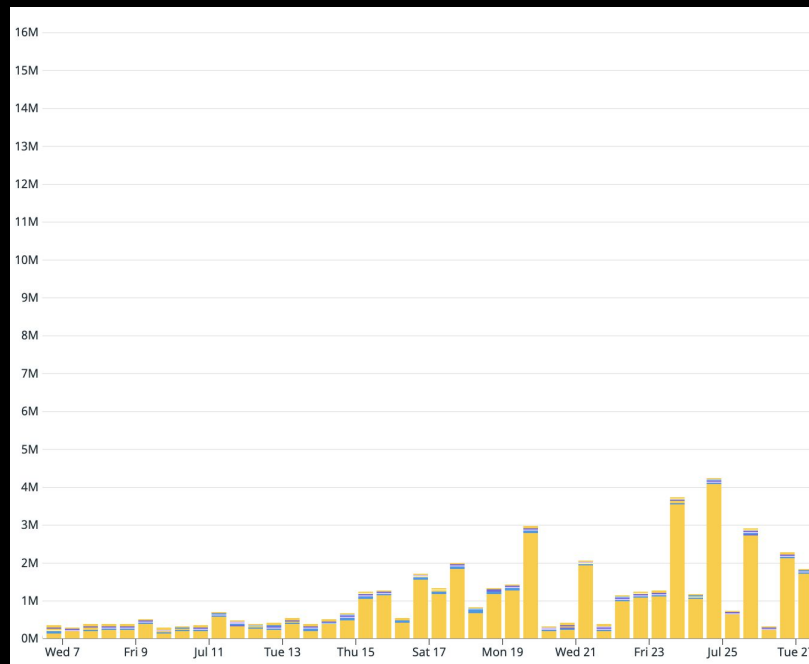
USE Method

Node with CPUs saturated & tcp retransmissions

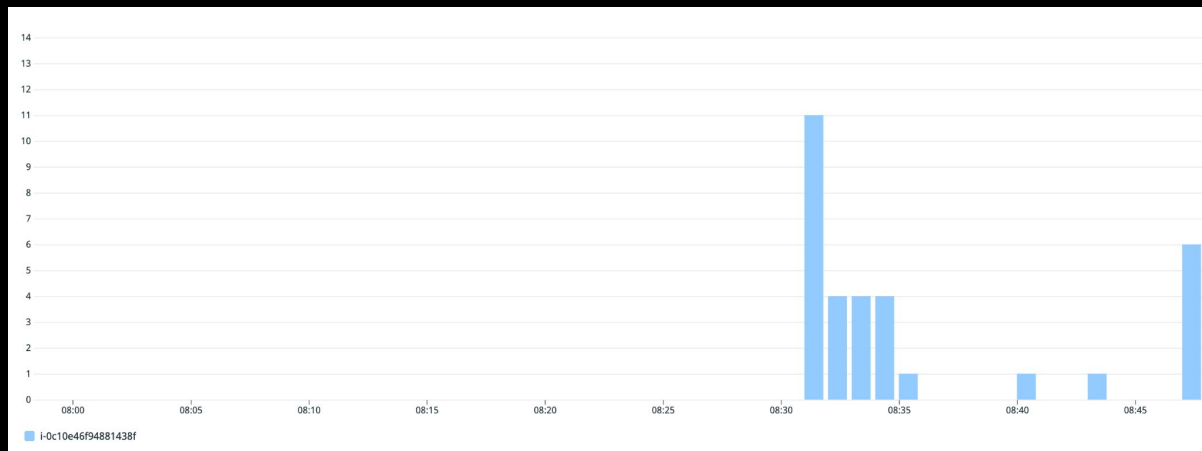


USE Method

Node with CPUs saturated & tcp retransmissions



I/O increases during the connectivity issue



Drill-Down

We can start from the software layer to the hardware one [Gregg 01].

We can split this analysis into three stages [McDougall 01]:

1. **Monitoring**
2. **Identification**
3. **Analysis**

Drill-Down

Five whys

Basically, we ask ourselves “why”, then answer the question and repeat it five times (or until finding the root cause).

1. Latency has increased on the cart endpoint. Why?
2. Lots of requests are queued in uWSGI. Why?
3. The horizontal pod autoscaler can't scale the deployment. Why?
4. Not enough VM instances to handle the load. Why?
5. The number of maximum VM is too low.

Drill-Down

Toolbox

Containerized environment... means restrictive
accesses... means harder to debug/access linux
primitives

BUT we have a tool named `ernestx` with
traditional tools and eBPF ones <3



Drill-Down

tcpretrans ebp output

Calico (179)

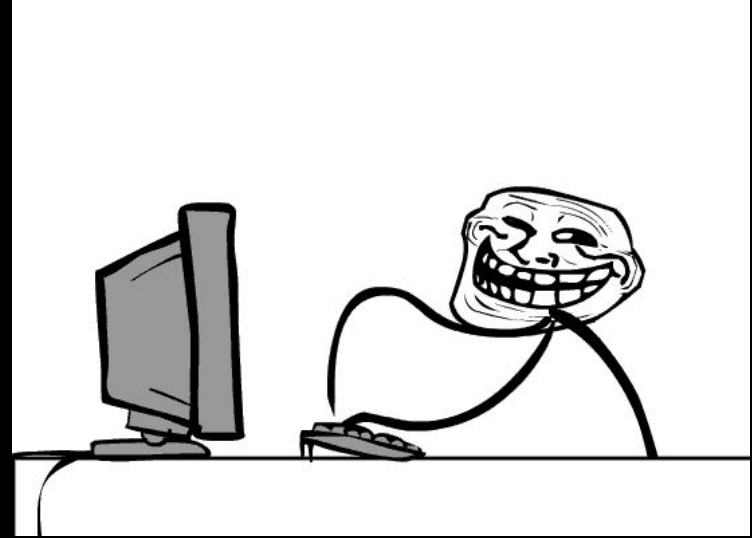
Redis (6379)

06:46:17 0		:51880		:6379	ESTABLISHED
06:46:19 0		:60153		:179	FIN_WAIT1
06:46:23 0		:56184		:6379	ESTABLISHED
06:46:37 0		:179		:35465	ESTABLISHED
06:46:37 0		:179		:35465	ESTABLISHED
06:46:38 8012		:179		:35465	ESTABLISHED
06:46:39 0		:179		:35465	ESTABLISHED
06:46:43 0		:179		:35465	ESTABLISHED
06:46:44 0		:37600		:6379	ESTABLISHED
06:46:49 0		:179		:35465	ESTABLISHED
06:47:03 0		:179		:35465	ESTABLISHED
06:47:10 0		:51880		:6379	ESTABLISHED
06:47:31 0		:179		:35465	ESTABLISHED

calico retrans

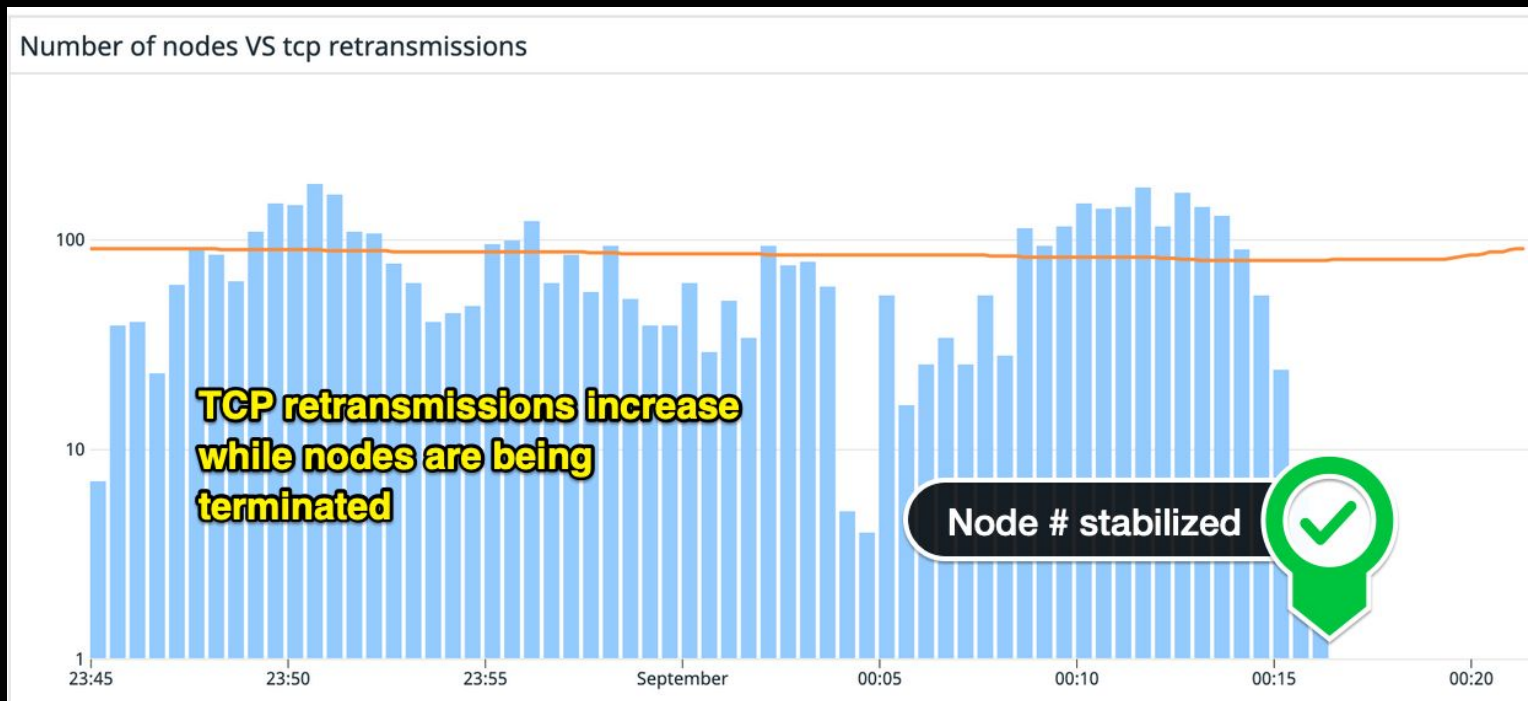
Why?

Calico uses BGP, you know the protocol that Facebook engineers doesn't know how to use it.



calico retrans

Why?



redis retrans

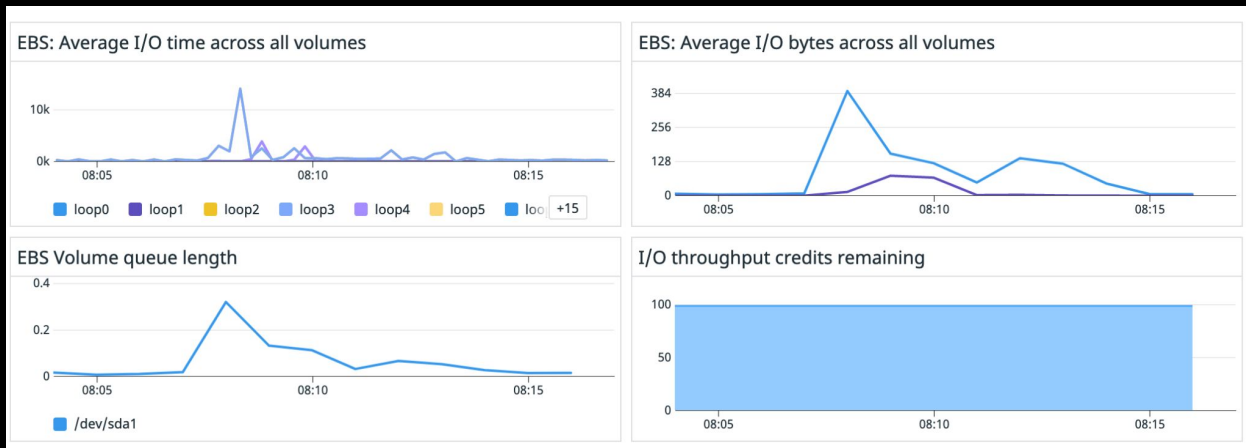
Why?

Disclaimer: no heavy disk I/O usage for
our nominal workload

redis retrans

Why?

Nothing obvious until we add a disks monitoring on our USE dashboard



i/o usage

Why?

Ernestx, my dear friend launch biosnoop at 8AM on an impacted instance

Here is a sample output of one of the biosnoop execution:

TIME(ms)	DISK	COMM	PID	LAT(ms)
1820	nvme0n1	agent	8926	26
9122	nvme0n1	dpkg	411280	24
9123	nvme0n1	dpkg	411280	24
9123	nvme0n1	dpkg	411280	24
9328	nvme0n1	dpkg	411280	25
9328	nvme0n1	dpkg	411280	25

dpkg

Why?

Tatu Ylönen, thank you for having invented the ssh protocol



Here is a sample output of the /var/log/dpkg.log file:

```
2021-08-24 06:49:41 startup packages configure
2021-08-24 06:49:41 configure libsystemd0:amd64 245.4-4ubuntu3.11 <none>
2021-08-24 06:49:41 status unpacked libsystemd0:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:41 status half-configured libsystemd0:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:41 status installed libsystemd0:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:41 startup packages configure
2021-08-24 06:49:41 configure systemd:amd64 245.4-4ubuntu3.11 <none>
2021-08-24 06:49:41 status unpacked systemd:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:41 status half-configured systemd:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:43 status installed systemd:amd64 245.4-4ubuntu3.11
2021-08-24 06:49:43 configure systemd-timesyncd:amd64 245.4-4ubuntu3.11 <none>
```

dpkg

Why?

```
$ systemctl list-timers
```

NEXT	LEFT	LAST	PASSED	UNIT
------	------	------	--------	------

ACTIVATES

Wed 2021-08-25 06:07:50 UTC 20h left	Tue 2021-08-24 06:30:42 UTC 2h 46min ago
--------------------------------------	--

apt-daily-upgrade.timer	apt-daily-upgrade.service
-------------------------	---------------------------

[....]

dpkg

Why?

[Timer]

OnCalendar=*-*-* 6:00

RandomizedDelaySec=60m



dpkg

Why?

- Node instance AMI was older than 6 months as we don't have **YET** a golden pipeline
- Happened each morning for the instances less older than 24h...



Wise words

& takeaways



BM Architecture

- Create a golden image with the greatest of care.
- I consider having knowledge on the eBPF eco-system is a must have in an SRE's toolbelt. We are performance investigators!
- Because entropy, randomness, disorder is part of human kind consider using methods while debugging a performance issue to avoid dispersing yourself.
- Write everything you are doing during a debugging session to avoid repeating yourself. **Think for the future you and your colleagues!!!**
- Why...why...why...why...why ; **always question everything.**

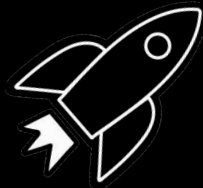
References

- [Gregg 01]: Gregg, B., *Systems Performance Second Edition*, Addison-Wesley, 2020
- [McDougall 01]: McDougall, R., Mauro, J., and Gregg, B., *Solaris Performance and Tools: DTrace and MDB Techniques for Solaris 10 and OpenSolaris*, Prentice Hall, 2006.
- https://en.wikipedia.org/wiki/Border_Gateway_Protocol
- <https://www.brendangregg.com/usemethod.html>



We are hiring

Chaos, SLO, Perf, perf, Perf, perf, Perf





Questions?

Remarks?

