

Linux Málaga

@linux_malaga www.linux-malaga.org



Bio Dala

cassandra

Bienvenido - Welcome - Witam

Juan Miguel Taboada Godoy

Juan José Soler Ruiz



@centrologic_es
http://linkedin.com/user/centrologic

@soleronline http://es.linkedin.com/in/soleronline





Juan Miguel Taboada Godoy (1980 - ...)

1996 – Primer ordenador y primera LAN (cable coaxial)

1999 – Universidad de Málaga y Linux Málaga

2001 – Grupo de investigación **GEB**.uma.es (4 años) Cluster computación +20 nodos (**OpenMosix**)

2002 – Presidente de Asociación **Málaga Wireless**

2003 – Beca en **Neurociencia** en SUNY Teleruta (Ministerio de Fomento – 2 años)

2005 – Autónomo:

- Nace Centrologic
- Polonia (2 años)
- Likindoy (Axaragua + **Junta Andalucía**)

2008 – Responsable Sistemas en PontGrup

2011 – Adquisición Datos en Bética Fotovoltáicas

2012 - SAFECLON y SCRUM/KANBAN

2013 – Executive MBA

2014 – Aeronáutica: Django + AngularJS

2015 – Industria: Likindoy + Big Data

Juan José Soler Ruiz

2001-2003 – CFGS Administración Sistemas Informáticos

2003 – Primer premio en el concurso "Javier Benjumea"

2003 – Montaje y configuración de "Cluster Heterogéneo De Computadoras" bajo SO Red Hat 7.2.

2005-2012 – STEA Telemática

2007-2009 – Primer CRM en PHP

2010-2011 – Administrador de BBDD / Programador Web en Bética Fotovoltáicas

2010-2012 – Opositometro

2012-....-Centrologic

2013 – Dailymarkets

2013-2014 - CRM en Python/Django

2014-... – Centrologic











Software de adquisición masiva de datos:

- 1 dato (Fecha+Valor) cada minuto por se al
- 1000 se ales por dispositivo20 dispositivos por nodo40 nodos por cliente

800.000 registros por minuto (8.000.000 tomas por minuto)

- 48M por hora1.152M por día

0,4 Billones por a Oycliente

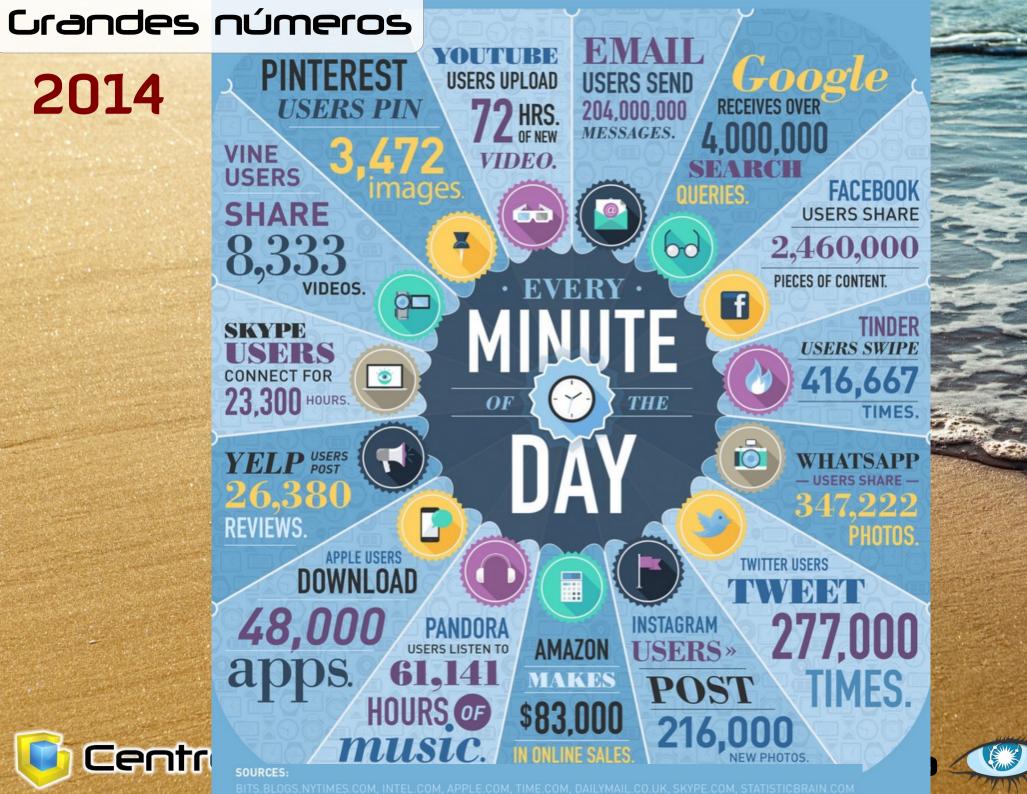








SOURCES: HTTP://NEWS.INVESTORS.COM/, ROYAL.PINGDOM.COM, BLOG.GROVO.COM Crandes números BLOG.HUBSPOT.COM, SIMPLYZESTY.COM, PCWORLD.COM, BIZTECHMAGAZINE.COM, DIGBY.COM 2012 MESSAGES. 2,000,000 SEARCH QUERIES. PIECES OF CONTENT. CONSUMERS ARE CREATED. FOURSQUARE USERS TWITTER USERS SEND OVER Centr



Ano 2016

- La gestión de datos crece de momento EXPONENCIALMENTE
- SEAGATE anuncia que el año 2016 será el año del Zettabyte
- 1 ZB equivale al espacio que ocupa 2 billones de años de música
- 1 ZB = 1024 exabytes = casi 1.1 trillones de Gbytes
- 1 zB = 1 sixtillón de bytes
- Ahora estamos en la época del **quintillón** de bytes...





¿Cuanto es un quintillón?

- Una única gota de agua contiene: 1.7 quintillones de moléculas de agua.
- La distancia de la Via Láctea hasta Andrómeda es de: 2 millones de años luz 18,87 quintillones de kilómetros 11,73 quintillon miles
- La tierra completa contiene unos: 1.234 quintillones de litros de agua 326 quintillon gallons of water
- Si cortamos la tierra por la mitad, la sección tendría un área aproximada de:
 - 1.275 quintillones de centímetros cuadrados
- ¿Cuanto es un quintillón de céntimos o peniques de dólar?



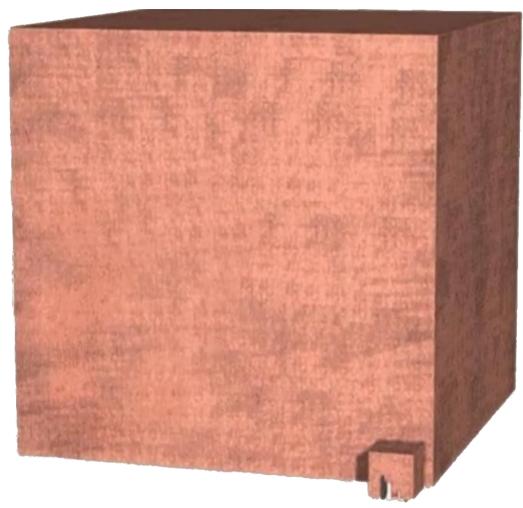


¿Cuanto es un quintillón... de peniques?





¿Cuanto es un quintillón... de peniques?



1.000.067.088.384.000.000 peniques 1 quintillón, 67 trillones, 88 billones, 384 millones de peniques Un cubo de **8,32 kilómetros** de lado





Teorema de CAP

Eric Brewer (2000)

Consistency Consistencia

Availability Disponibilidad A:

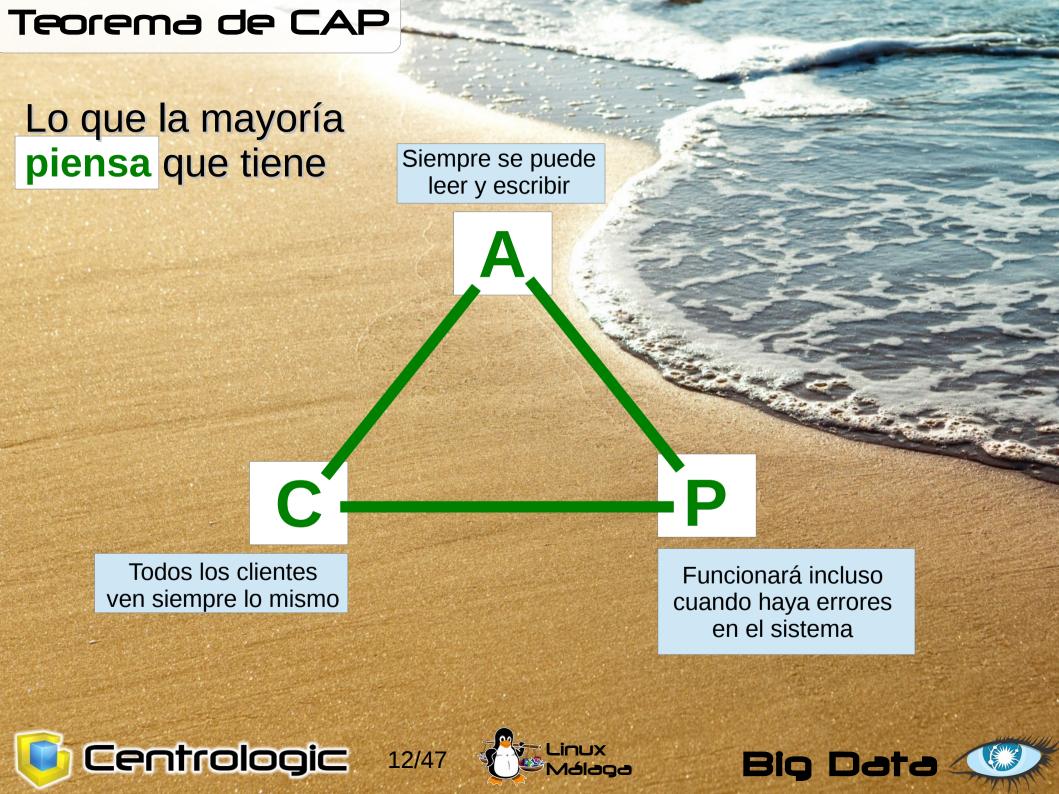
P: Partition tolerance → Tolerancia al particionado

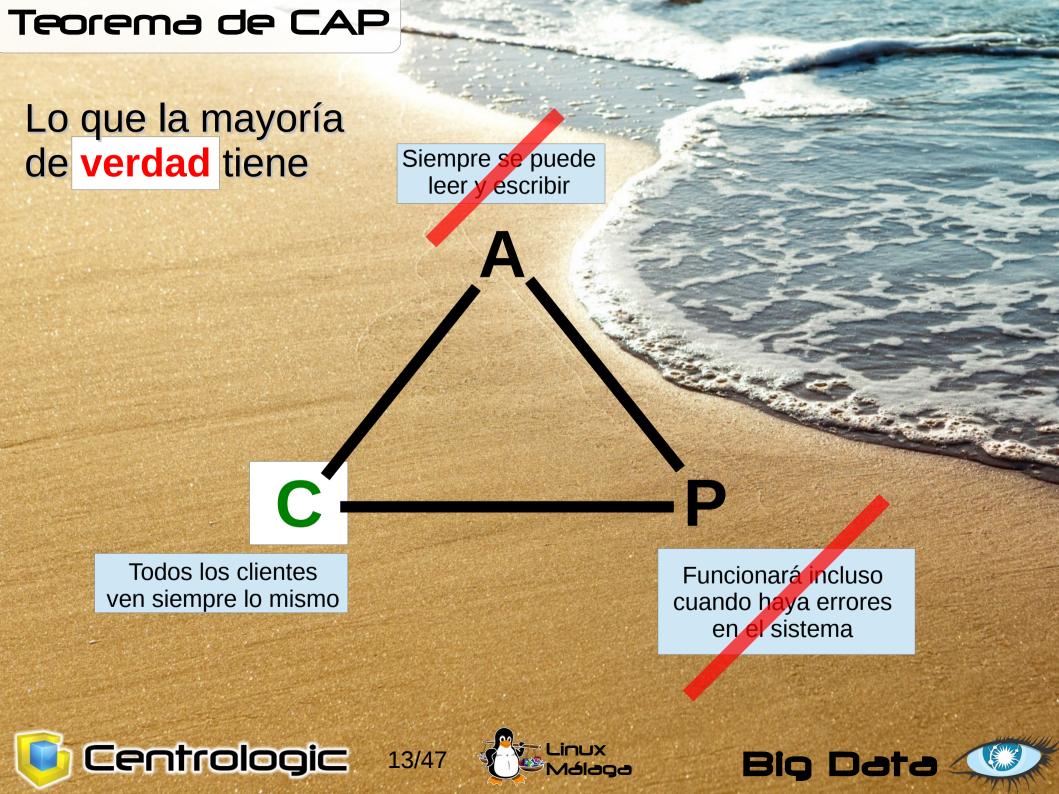
Sólo puedes llegar a 2 de las 3

Es imposible para un sistema de cómputo distribuido garantizar simultáneamente la consistencia, la disponibilidad y ser tolerante al particionado de los datos (separación y distribución).









Teorema de CAP

Lo cierto es que todos buscamos la disponibilidad (A)

Pero ... ¡¡¡ tenemos que elegir entre... !!!

Escalabilidad (P)

y

Consistencia (C)







ACID

A: Atomicidad

C: Consistencia

I: Aislamiento (Isolation)

D: Durabilidad

En grandes sistema ocurre que: Disponibilidad y Rendimiento





BASE

BA: Básicamente disponible

S: Flexible (Soft state)

E: Consistencia eventual

Da menos importancia a la consistencia en pro de la tolerancia al particionado aparece la consistencia eventual





¿Qué es la consistencia eventual?

Que ... eventualmente será consistente

Podemos introducir un dato y que no esté disponible inmediatamente después

Convergencia natural a la consistencia





Resolución de conflictos

Anti-entropía (control de versiones)

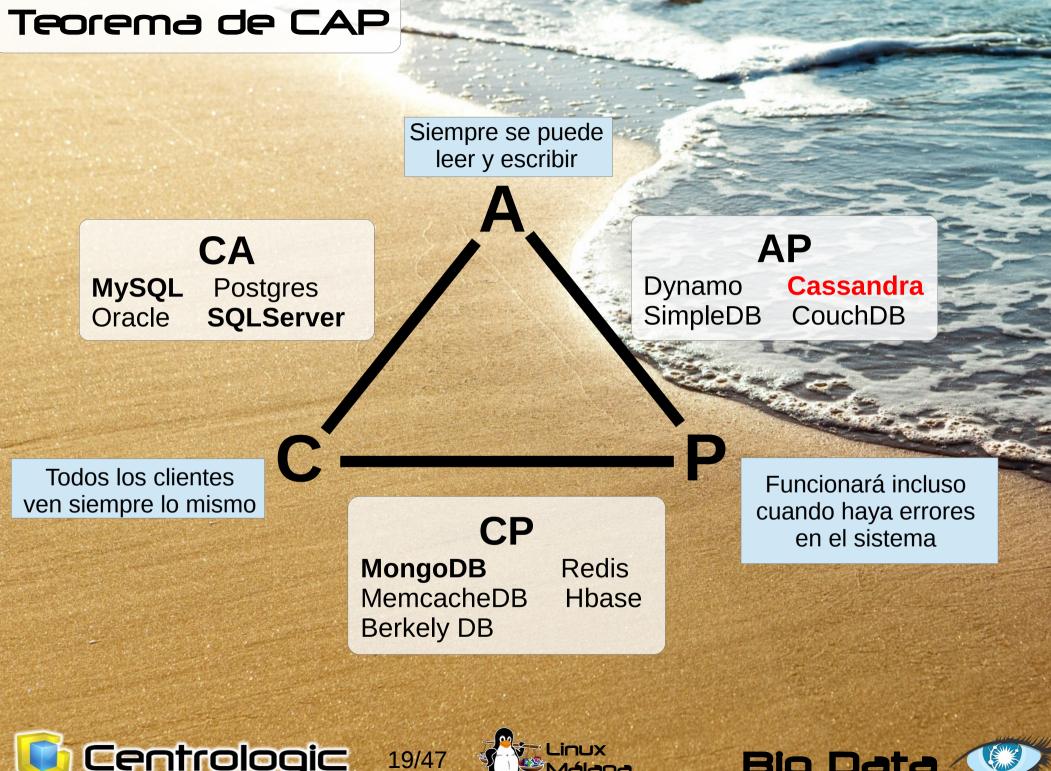
Reconciliación (elección de estado final) [generalmente: "last write wins"]

Garantizar la seguridad de una operación

Strong eventual consistency (SEC)













Teorema de CAP CA AP Cassandra **MySQL Dynamo Aster Data Voldemort SimpleDB Postgres** Greenplum **Oracle Tokyo Cabinet** CouchDB **Vertica SQLServer** Riak KAI **Tipo** Relacional Clave-Valor Orientado a columna Orientado a documento CP **BigTable Berkely DB MongoDB MemcacheDB Hypertable Terrastore Scalaris HBase** Redis











Disponibilidad contínua

Simplicidad en la gestión entre servidores Sin un único punto de fallo

Escalabilidad linear

- Si 2 nodos procesan 100 transacciones/seg
- 4 nodos procesan 200 transacciones/seg
- 8 nodos procesan 400 transacciones/seg

Sistema descentralizado (sin master)

Relaciones por grupo: Nodo → Datacenter → Cluster Replicación personalizada









Gossip: mantiene la red informada

Partitioner: deciden como se distribuyen los datos

Replication factor: número de réplicas en el cluster

Replica placement strategy: donde poner las réplicas

Snitch: define grupo de máquinas destinadas a réplicas







Gossip:

- Protocolo punto a punto
- Comunicación entre nodos
- Detección de fallos y recuperación
- Autodetección de topología e información







Particionador: encargado de "esparcir" los datos.

Desaconsejados:

Random: Hashes con MD5 (mejor usar Murmur3) ByteOrdered: orden lexicográfico sobre las claves OrderPreserving: se asumen claves en formato UTF8

Murmur3: Funcionalmente idéntico a Random pero es más rápido sin efectos colaterales. Se centra en la distribucción espacial.







Planificar el deploy de un cluster

Hardware:

- RAM: 16GB-64Gb (mínimo 8Gb)
- CPU: 8-cores dedicados o 4-8 cores en virtuales
- Disco: mejor entre 500Gb y 1Tb por nodo (según I/O)
- 2 discos (commit log + data)
- Sistema de ficheros XFS
- Red mínimo Gigabit









Planificar el deploy de un cluster

Espacio útil disco: 45%-70% de espacio total RAW

Antipatrones:

- Usar un NAS
- Sistemas de ficheros compartidos
- SELECT ... IN
- Leer antes de escribir (multiples hits)
- Balanceadores de carga
- Falta de testing
- Bajo conocimiento de Linux









SQL

USE myDatabase;

/* Creating Tables */ CREATE TABLE IF NOT EXISTS myTable (id INT PRIMARY KEY);

/* Altering Tables /* ALTER TABLE myTable ADD myField INT;

/* Creating Indexes */ CREATE INDEX myIndex ON myTable (myField);

/* Inserting Data */ INSERT INTO myTable (id, myField) VALUES (1, 7);

/* Selecting Data */ SELECT * FROM myTable WHERE myField = 7;

/* Counting Data */ SELECT COUNT(*) FROM myTable;

/* Deleting Data */ DELETE FROM myTable WHERE myField = 7;

COL

USE myDatabase;

/* Creating Tables */ CREATE TABLE IF NOT EXISTS myTable (id INT PRIMARY KEY);

/* Altering Tables /* ALTER TABLE myTable ADD myField INT;

/* Creating Indexes */ CREATE INDEX myIndex ON myTable (myField);

/* Inserting Data */ INSERT INTO myTable (id, myField) VALUES (1, 7);

/* Selecting Data */ SELECT * FROM myTable WHERE myField = 7;

/* Counting Data */ SELECT COUNT(*) FROM myTable;

/* Deleting Data */ DELETE FROM myTable WHERE myField = 7;











USE miBaseDatos;

/* Creando Tablas */ CREATE TABLE IF NOT EXISTS miTabla (id INT PRIMARY KEY);

/* Cambiando Tablas /* **ALTER TABLE miTabla ADD miCampo INT;**

/* Creating Indexes */ **CREATE INDEX milndice ON miTabla (miCampo)**;









```
/* Insertando Datos */
INSERT INTO miTabla (id, miCampo) VALUES (1, 7);
```

/* Seleccionando Datos */ **SELECT * FROM miTabla WHERE miCampo = 7;**

/* Contando Datos */ **SELECT COUNT(*) FROM miTabla;**

/* Borrando Datos */ **DELETE FROM miTabla WHERE miCampo = 7**;









/* Crear un nuevo keyspace en CQL */ **CREATE KEYSPACE miBaseDatos WITH replication =** {'class': 'SimpleStrategy', 'replication_factor': 1};

/* Crear una nueva base de datos en SQL */ **CREATE DATABASE miBaseDatos**;

No existen:

JOIN, GROUP BY, y FOREIGN KEY







- 1.- Las escrituras son baratas. Escribe todo del modo en que vas a leerlo.
- 2.- UPSERT: INSERT (inserta ó actualiza), UPDATE (inserta ó actualiza), ya que Cassandra no hace un read durante estos procesos.
- 3.- Filas con TTL, pasado X tiempo la fila caduca.
- 4.- DELETE ... ¡no borra!







/* Seleccionar datos de un rango */ **SELECT * FROM myTable** WHERE miCampo > 5000 AND miCampo < 100000;

Bad Request: Cannot execute this query as it might involve data filtering and thus may have unpredictable performance. If you want to execute this query despite the performance unpredictability, use ALLOW FILTERING.





¿Cómo se gestiona esto?

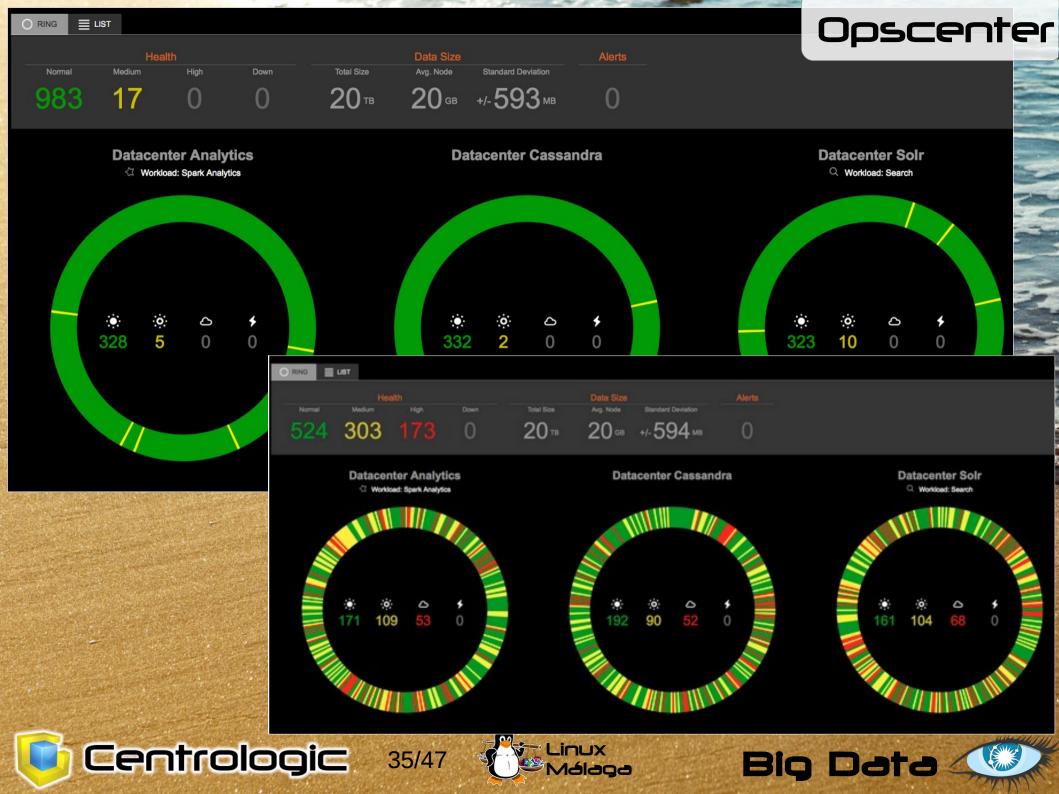


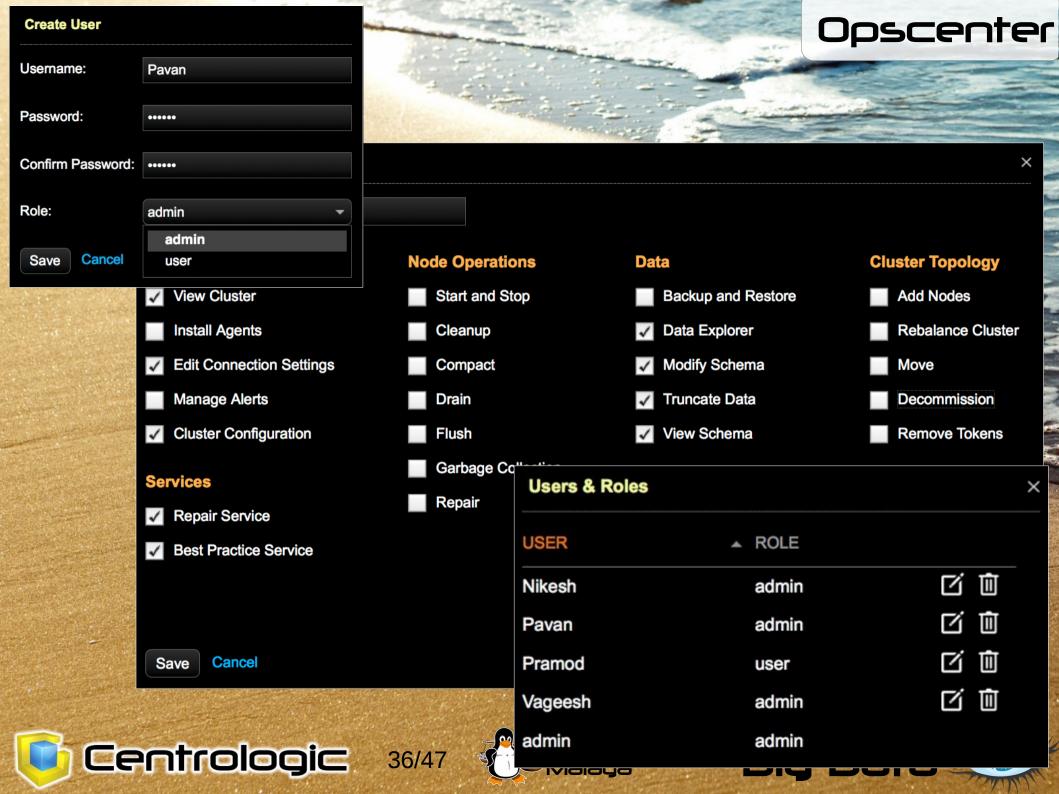












Opscenter

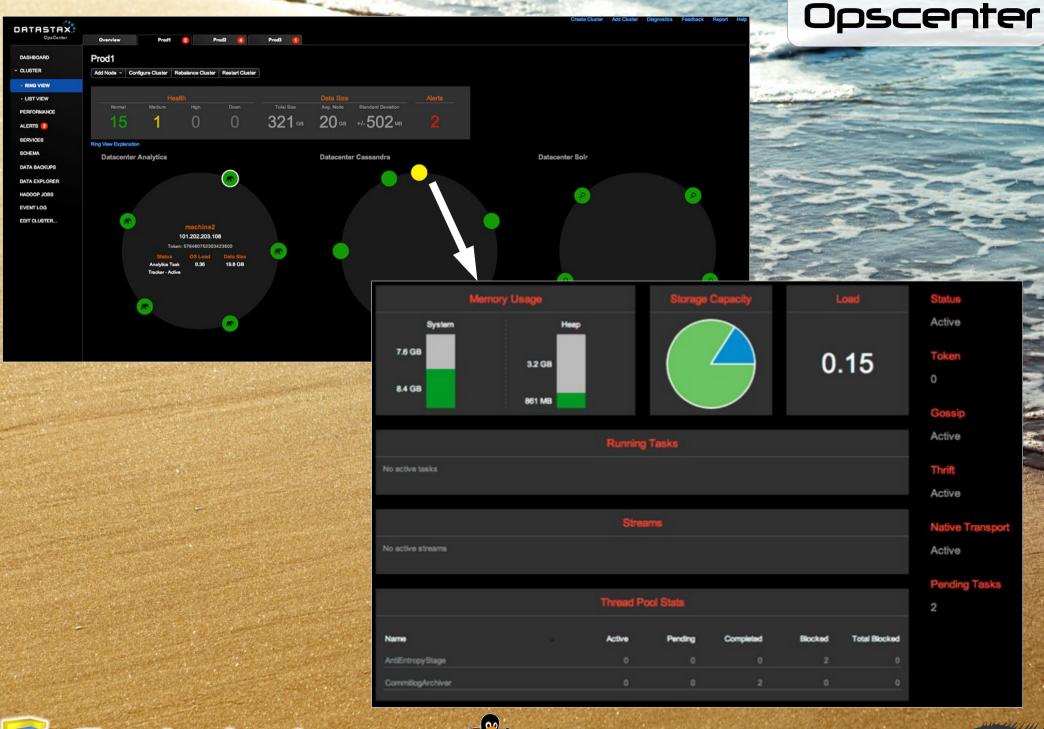










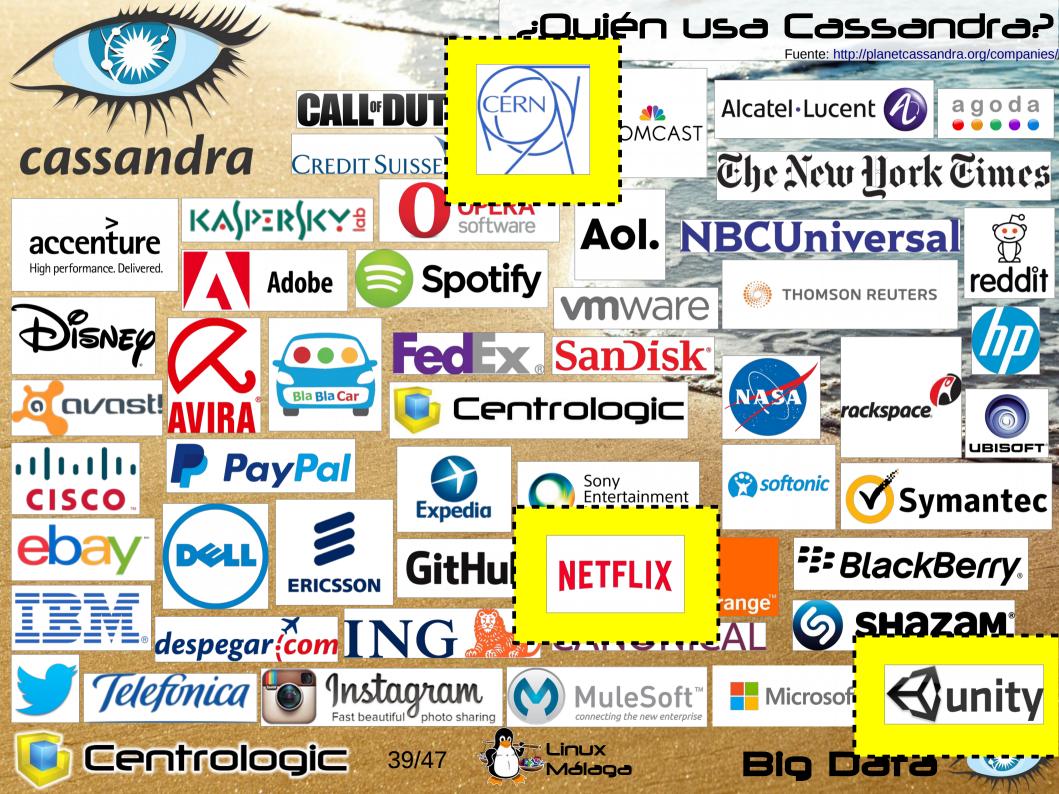














¿Quién usa Cassandra?





Al finalizar la última iteración del LHC en el que se descubrió el "Bosón de Higgs", el CERN almacenaba más de 100Pbytes

El LHC del CERN arrancó en Abril de 2015 en busca de la Supersimetría (capaz de producir 1Pbytes/segundo)

1Petabyte = 1.000 Terabytes = 1 Millón de Gigabytes

NETFLIX

Ejecuta 235 clusters separados, con un total de 7.000 nodos. 1 Millón de escrituras por segundo (factor 3)



Cunity Ecosistema para desarrolladores de juegos que teniendo problemas con MongoDB migraron a Cassandra.







...esta historia parecía de cuento...

...los datos se repartían por los nodos...

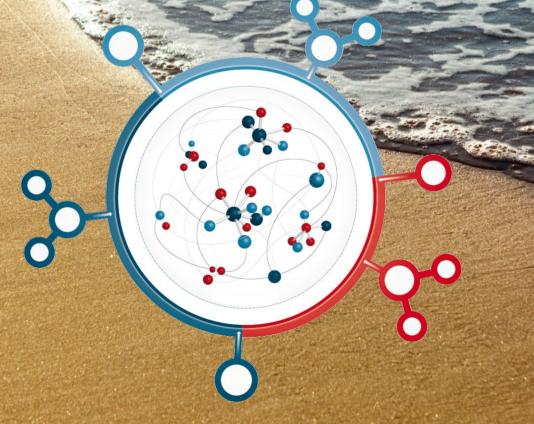
...existía resilencia...

...había alta disponibilidad...

...crecía linealmente f(nodos)...

...todo parecía ideal...











...pero
con el paso del tiempo
a la bella Cassandra
le surgió un amante...









...con el paso del tiempo ...conoció a KairosDB















"No podemos confundir potencia con potencial"



- Brian Hawkins (Salt Lake) y Jeff Sabin (Utah) de Proofpoint. (Link)
- Base de datos para series temporales sobre Cassandra
- 11 de abril de 2013 primera versión Beta
- Cada registro tiene: Métrica + Fecha + Tipo + Tags
- 3 semanas de datos o 1.814,4 Millones de columnas exactamente
- Genial...pero ¿cómo mejora esto el rendimiento?













WHEN IN DOUBT TRY ANOTHER HOLE









































Muchas Cracias



Juan Miguel Taboada Codoy http://www.centrologic.com

@centrologic es

http://linkedin.com/user/centrologic



Juan José Soler Ruiz

@soleronline

http://es.linkedin.com/in/soleronline



