

An underwater photograph showing a diver on the right side, wearing a yellow tank and black gear, swimming in clear blue water. A large, thick tree branch with green leaves extends from the left side across the middle of the frame. The water is bright and clear, with some bubbles visible near the diver.

20.000 reqs de viaje submarino

Bienvenido - Welcome - Witam



Juanmi Taboada
@juanmitaboada
www.juanmitaboada.com

OpenSouthCode'19 - Málaga, 25 may 2019

1996



1997



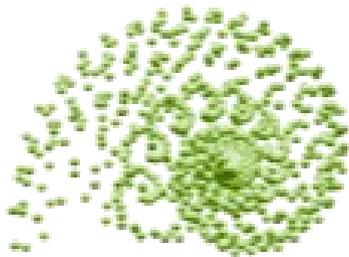
1999



UNIVERSIDAD
DE MÁLAGA

| uma.es

1999



grupo de estudios
en biomimética

2001



2005

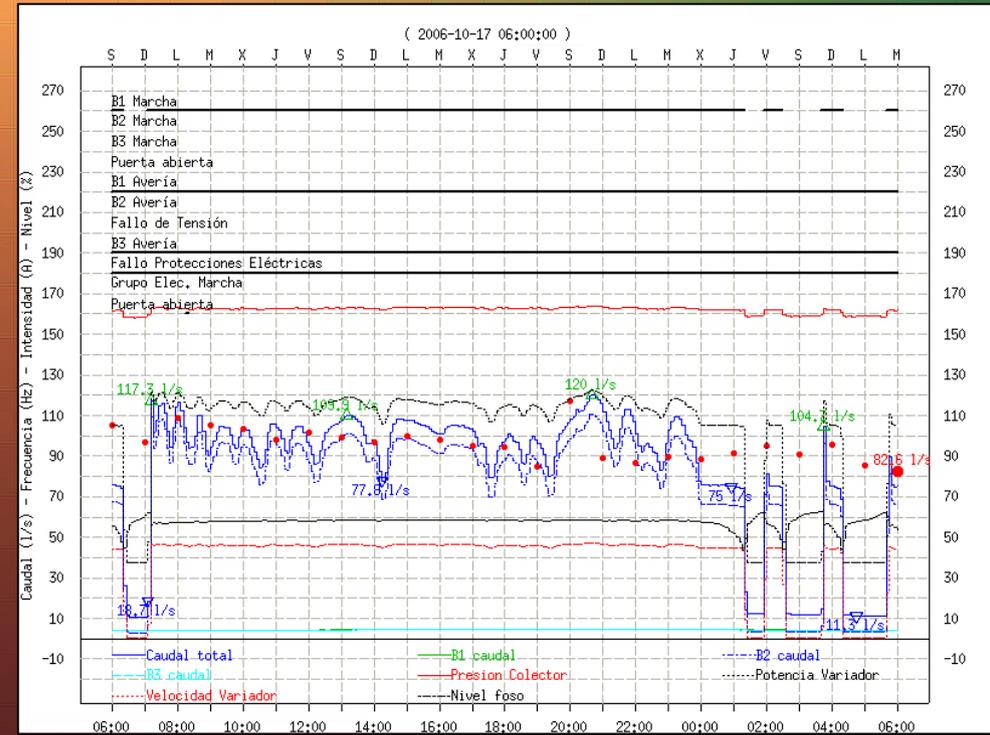
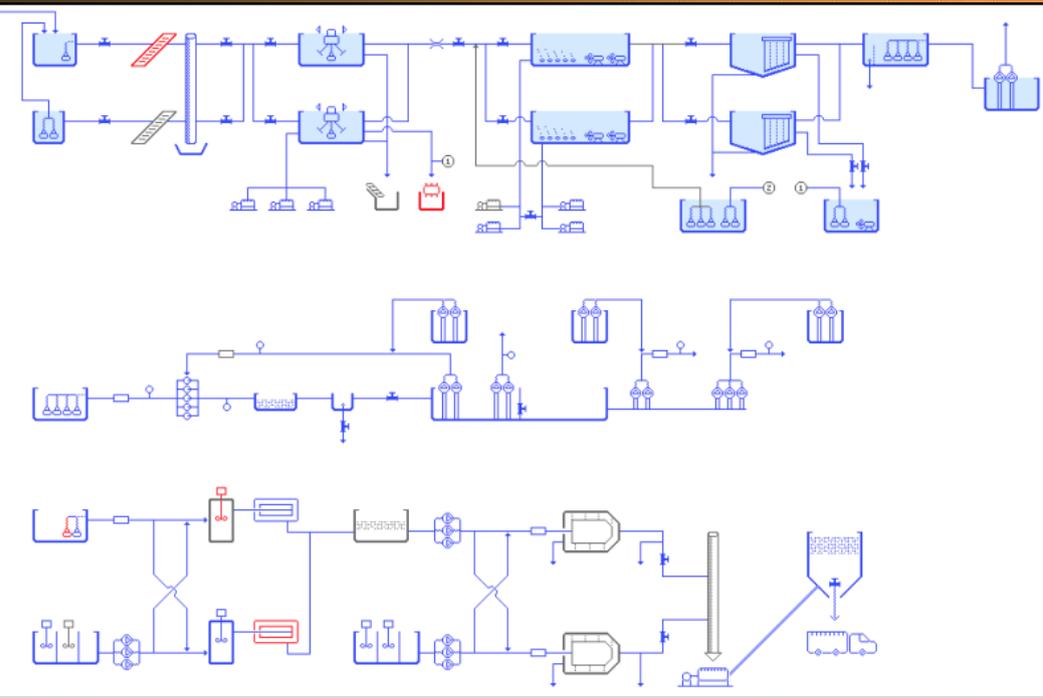


python



Likindoy

The monitoring company



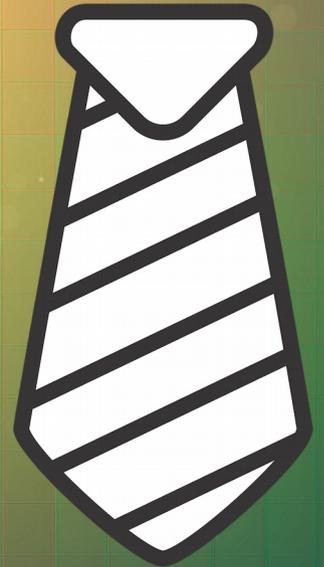


2008

django



Seguros
Webservices



EXECUTIVE
MBA



Fotovoltaica
Adquisición Datos

Aeronáutica



- Hoy -

20.000 regs/seg

o más ... porque no los contamos





¿Información?

Información

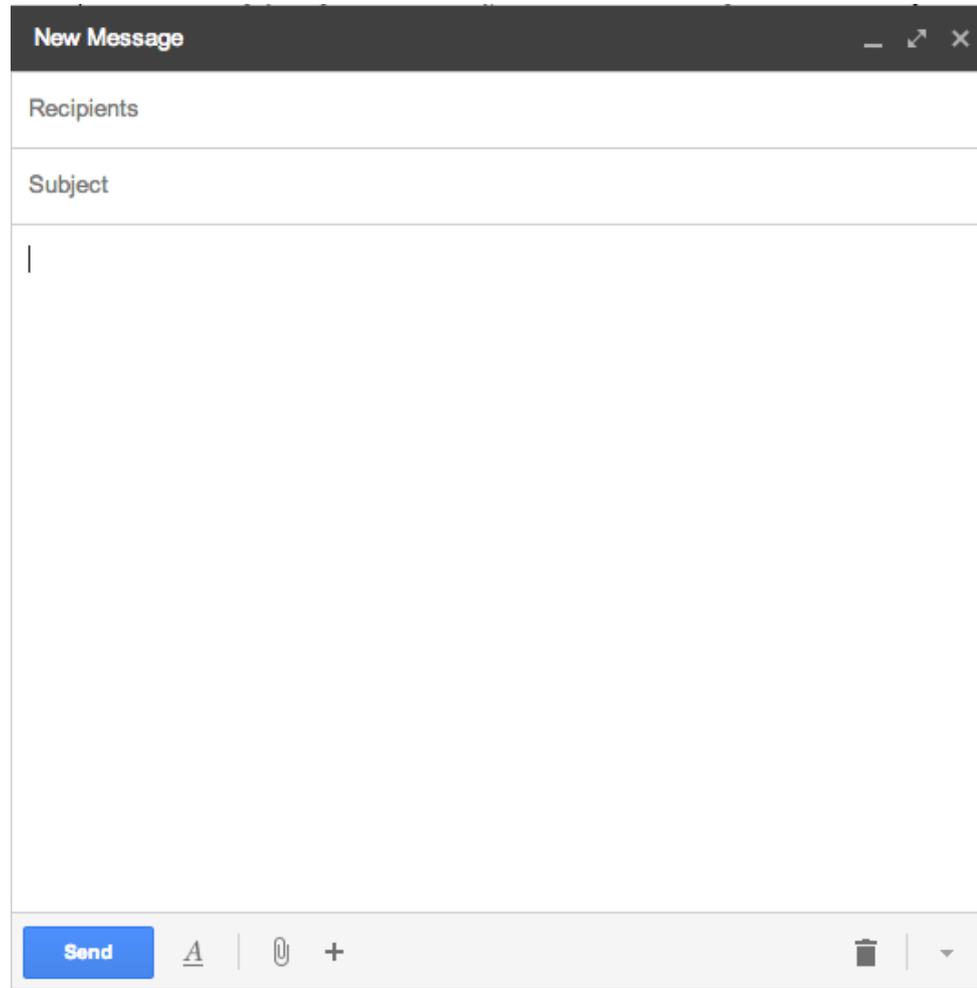
es un conjunto organizado de datos procesados,
que constituyen un mensaje
que cambia el estado de conocimiento
del sujeto o sistema que recibe dicho mensaje.

(fuente Wikipedia)

Información digital

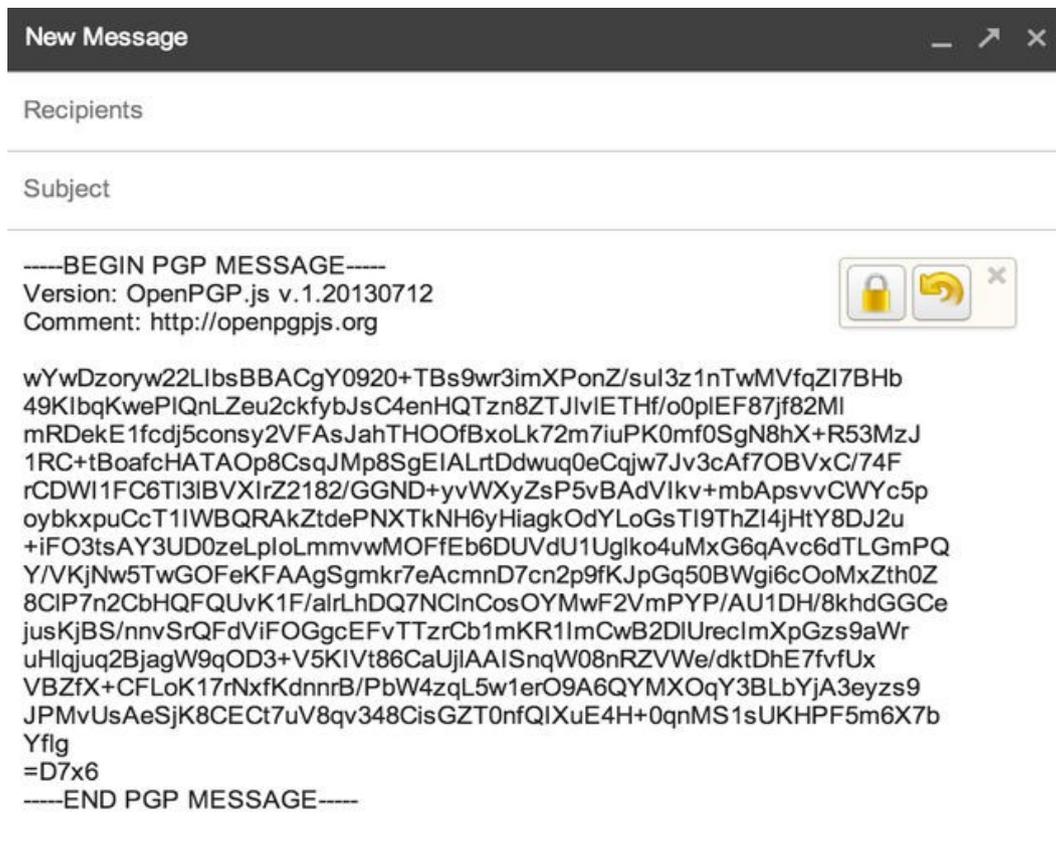


Información digital → Email

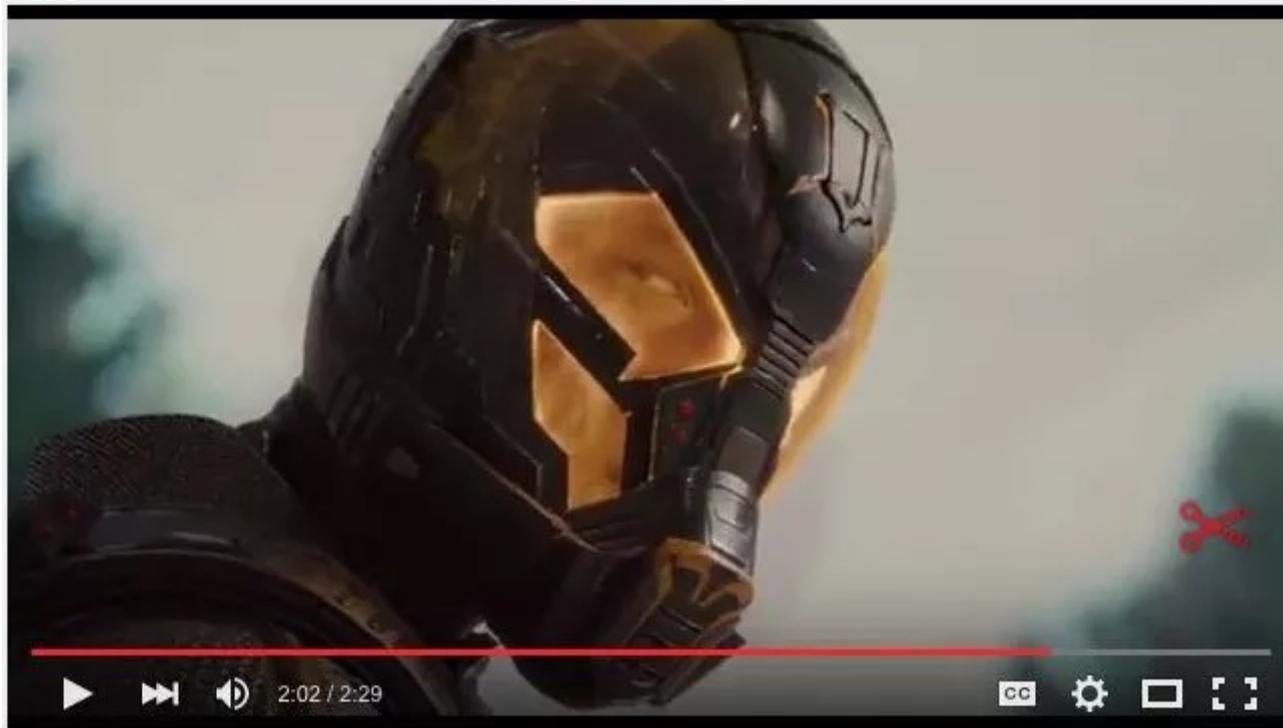
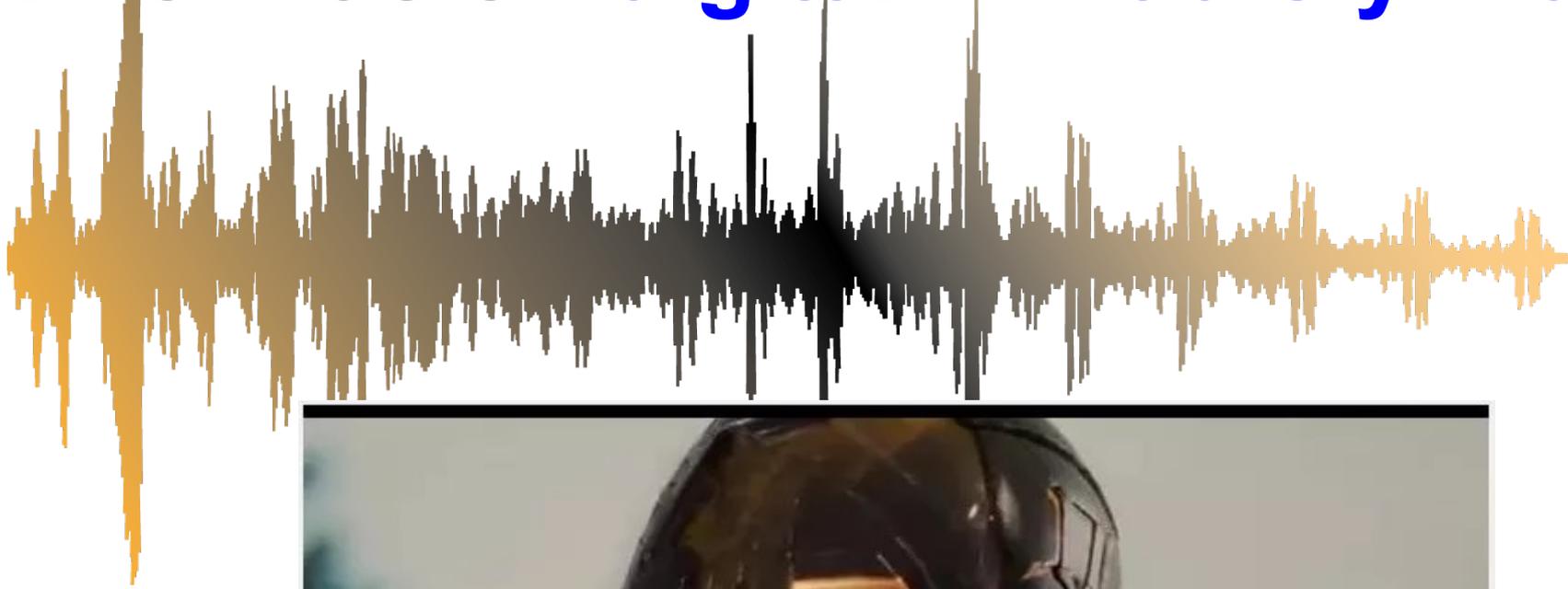


A screenshot of a 'New Message' email composition window. The window has a dark title bar with the text 'New Message' and standard window controls (minimize, maximize, close). Below the title bar, there are three main sections: 'Recipients', 'Subject', and a large text area for the message body. The 'Recipients' and 'Subject' fields are currently empty. The text area contains a single vertical cursor. At the bottom of the window, there is a toolbar with a blue 'Send' button, a text formatting icon (underline), an attachment icon, a plus sign, a trash icon, and a dropdown arrow.

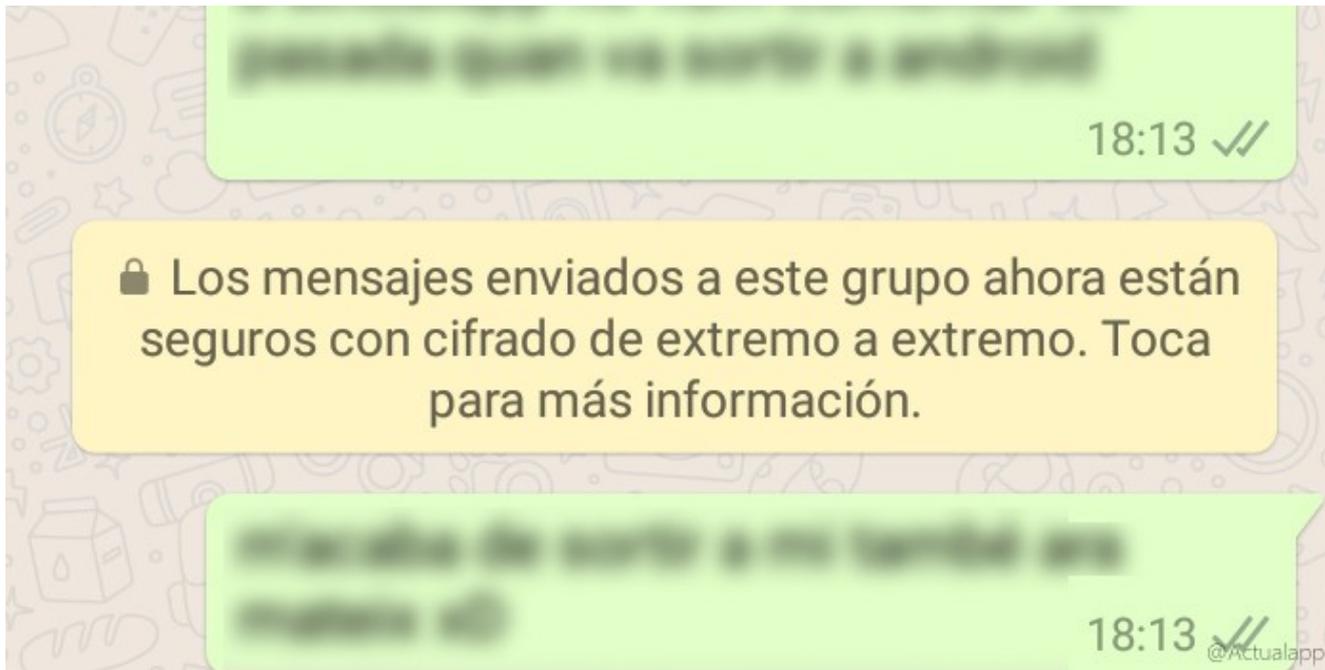
Información digital → Email



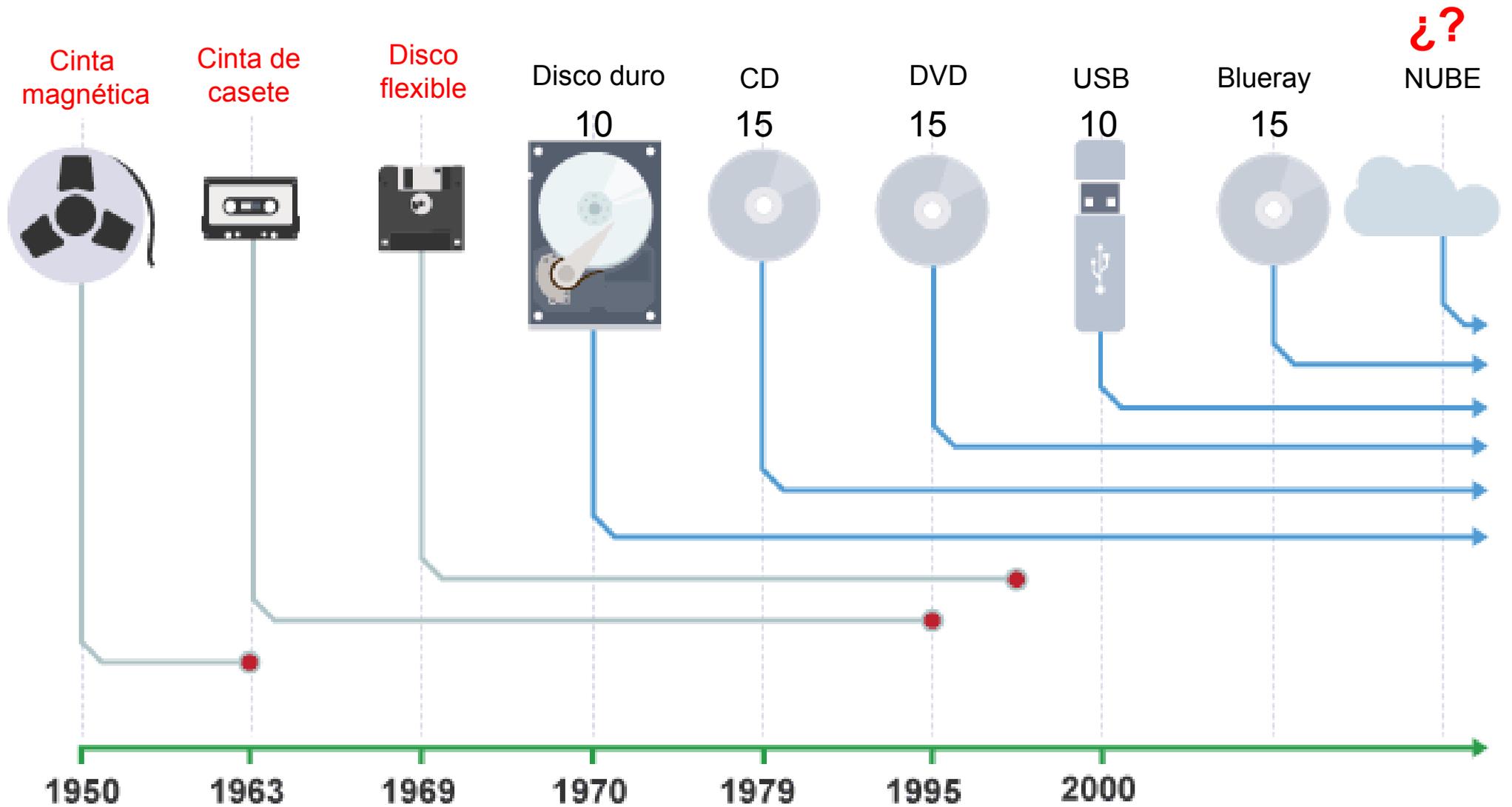
Información digital → Audio y Video



Información digital → Email



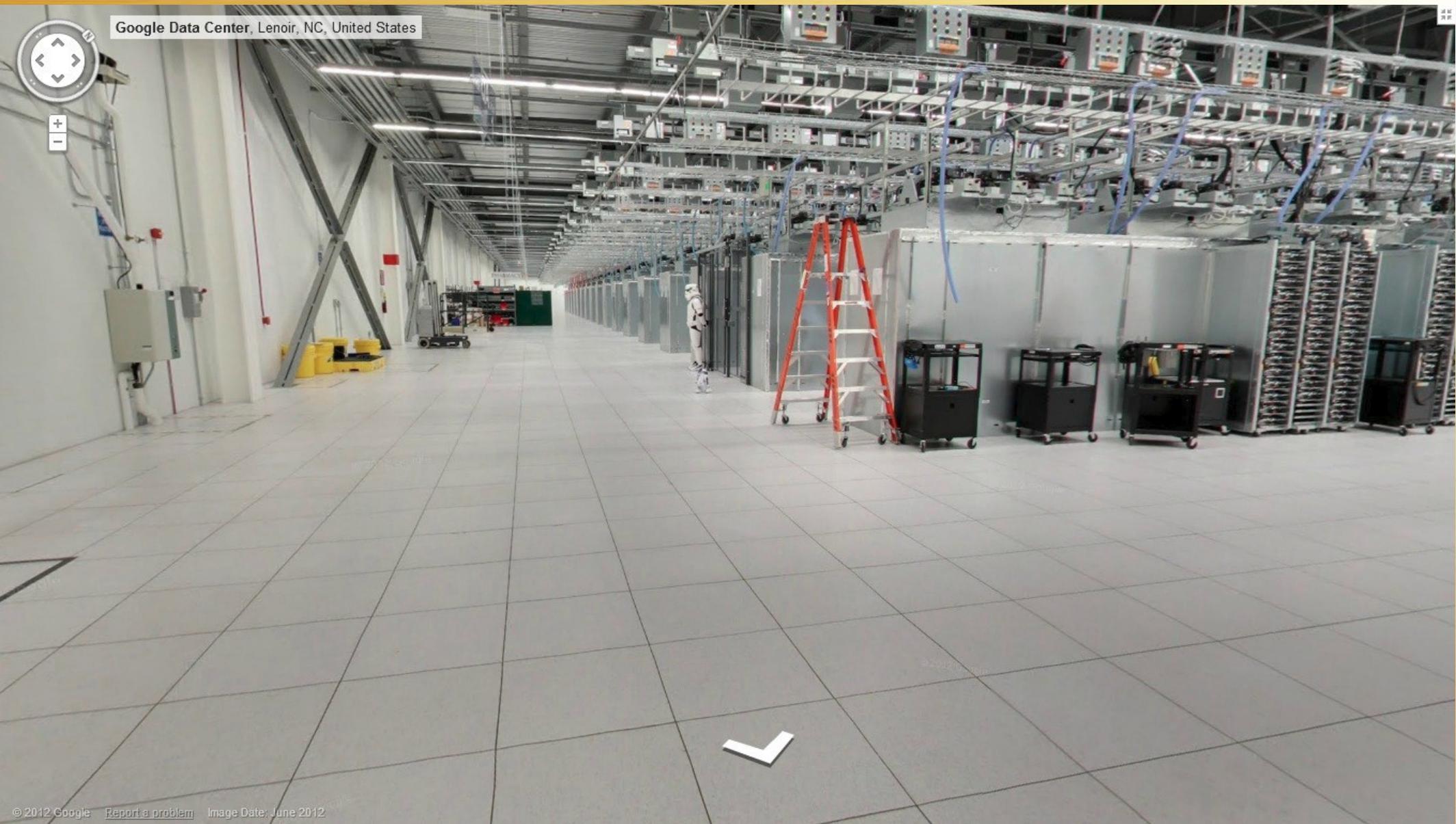
¿Información?



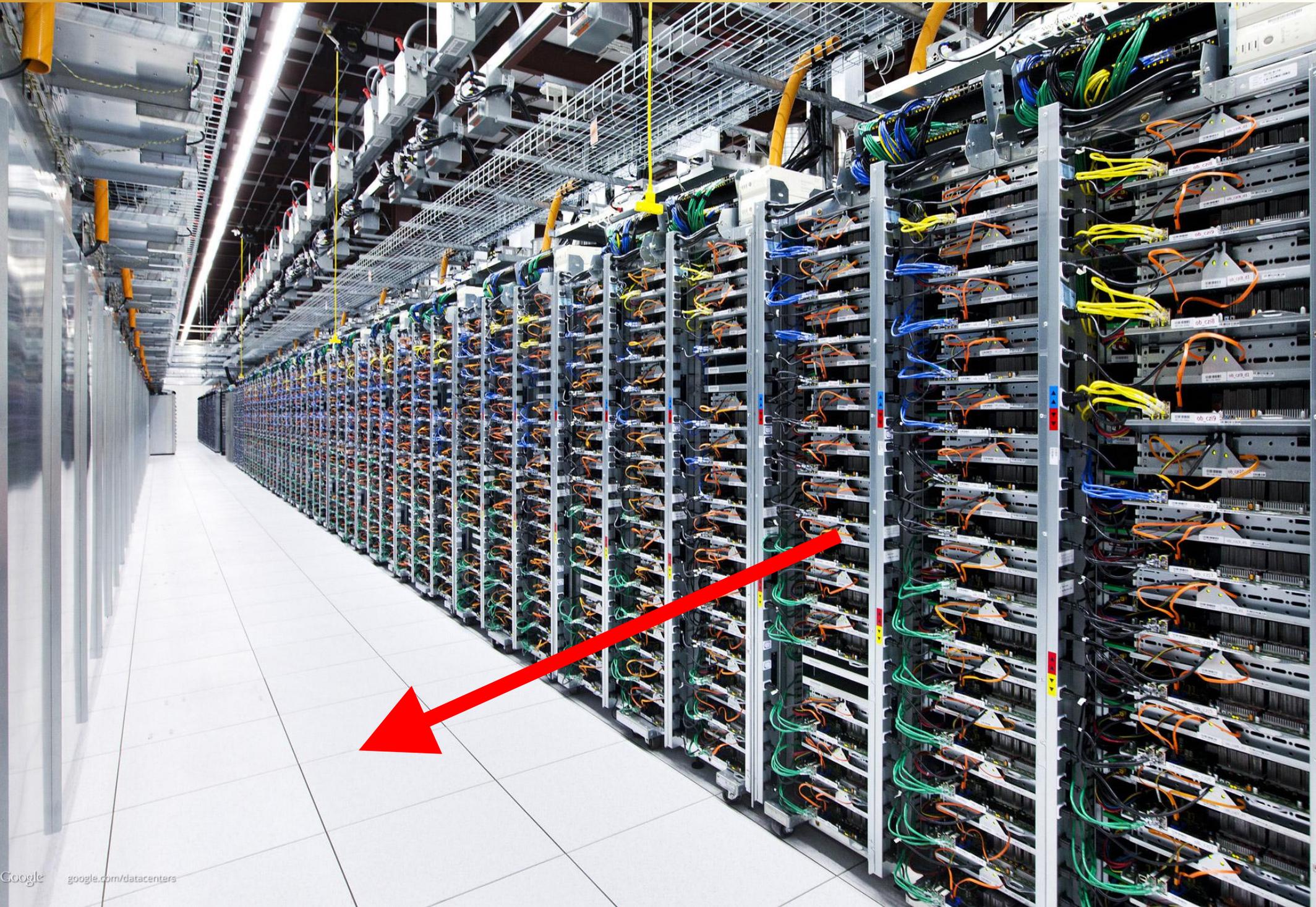
La NUBE



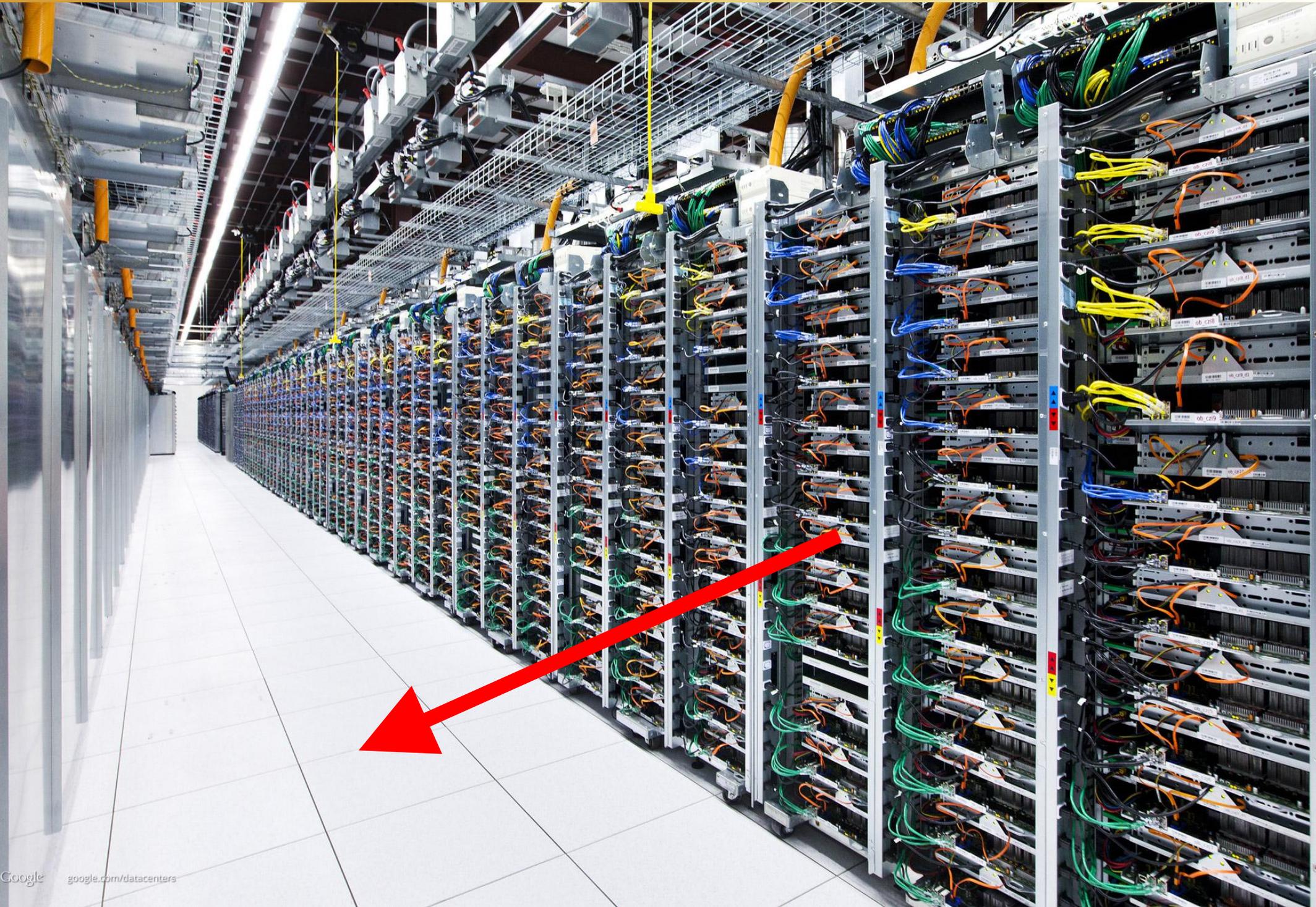
La NUBE: Datacenter



La NUBE: saquemos un equipo



La NUBE: saquemos un equipo

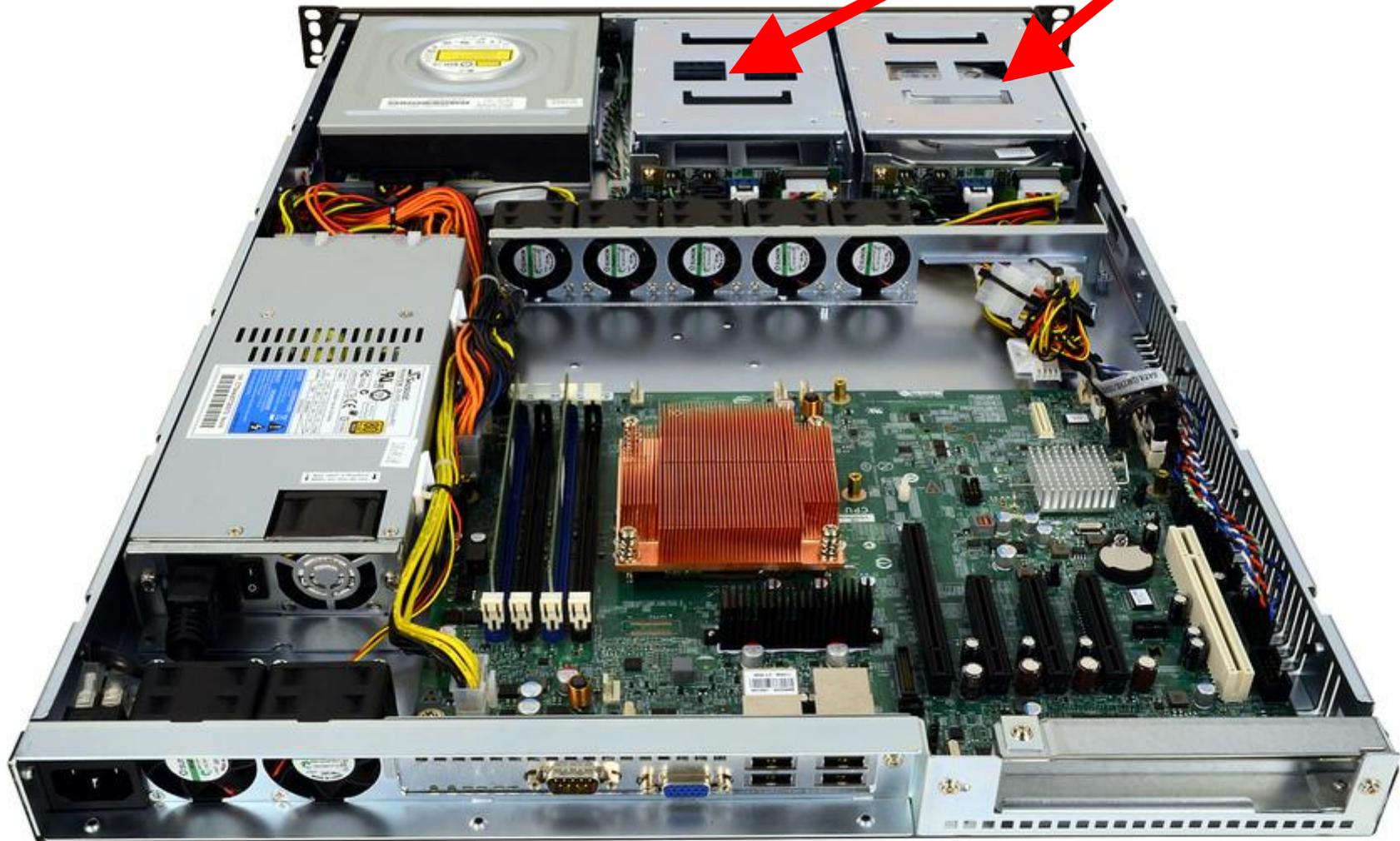


La NUBE: abramos el equipo



La NUBE: interior del ordenador

Disco duro



1 disco duro = 10 años = 3 TBytes



10 años



¡Sincronicemos todo el tiempo!



¿Cómo pinta todo esto?

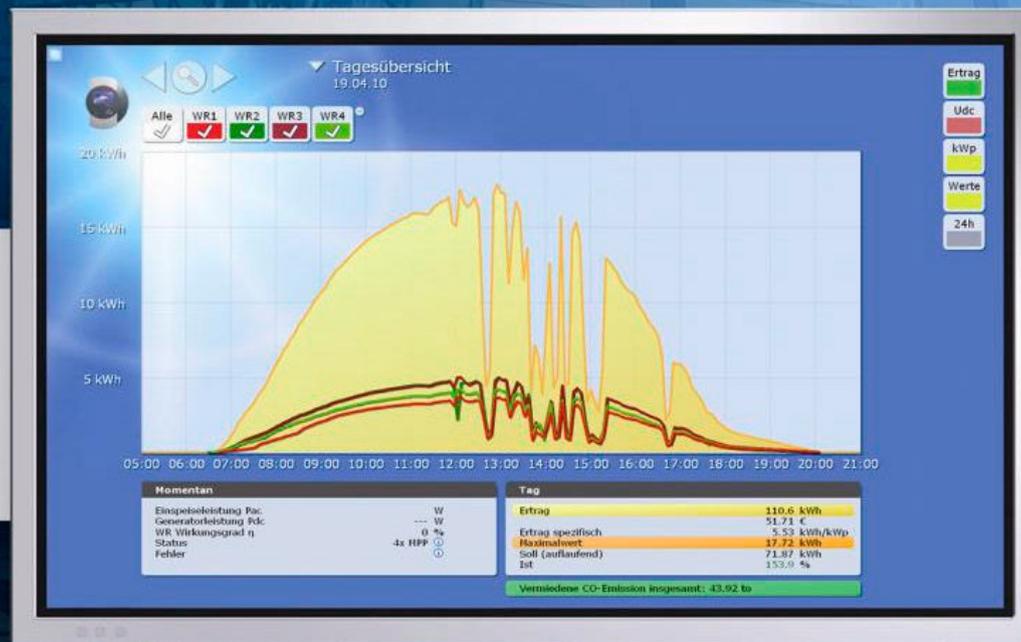


Juanmi Taboada

Big Data



Monitorización



Plant example 2

Zeit	Typ	Wärmeleistung	WR
18.01.2010 22:30:16	Verbindung	Wärme Daten vom WR 1 WR Typ: WRVCL16 WR Abk: 2000270907	WRVCL16 2000270907
17.01.2010 22:30:36	Verbindung	Wärme Daten vom WR 1 WR Typ: WRVCL16 WR Abk: 2000270907	WRVCL16 2000270907
16.01.2010 22:30:16	Verbindung	Wärme Daten vom WR 1 WR Typ: WRVCL16 WR Abk: 2000270907	WRVCL16 2000270907
15.01.2010 22:30:36	Verbindung	Wärme Daten vom WR 1 WR Typ: WRVCL16 WR Abk: 2000270907	WRVCL16 2000270907

La cruda realidad

Recibimos :

- 1 registro tiene 100 datos útiles (18 campos)
- 20.000 registros por segundo

1.200.000 registros por minuto

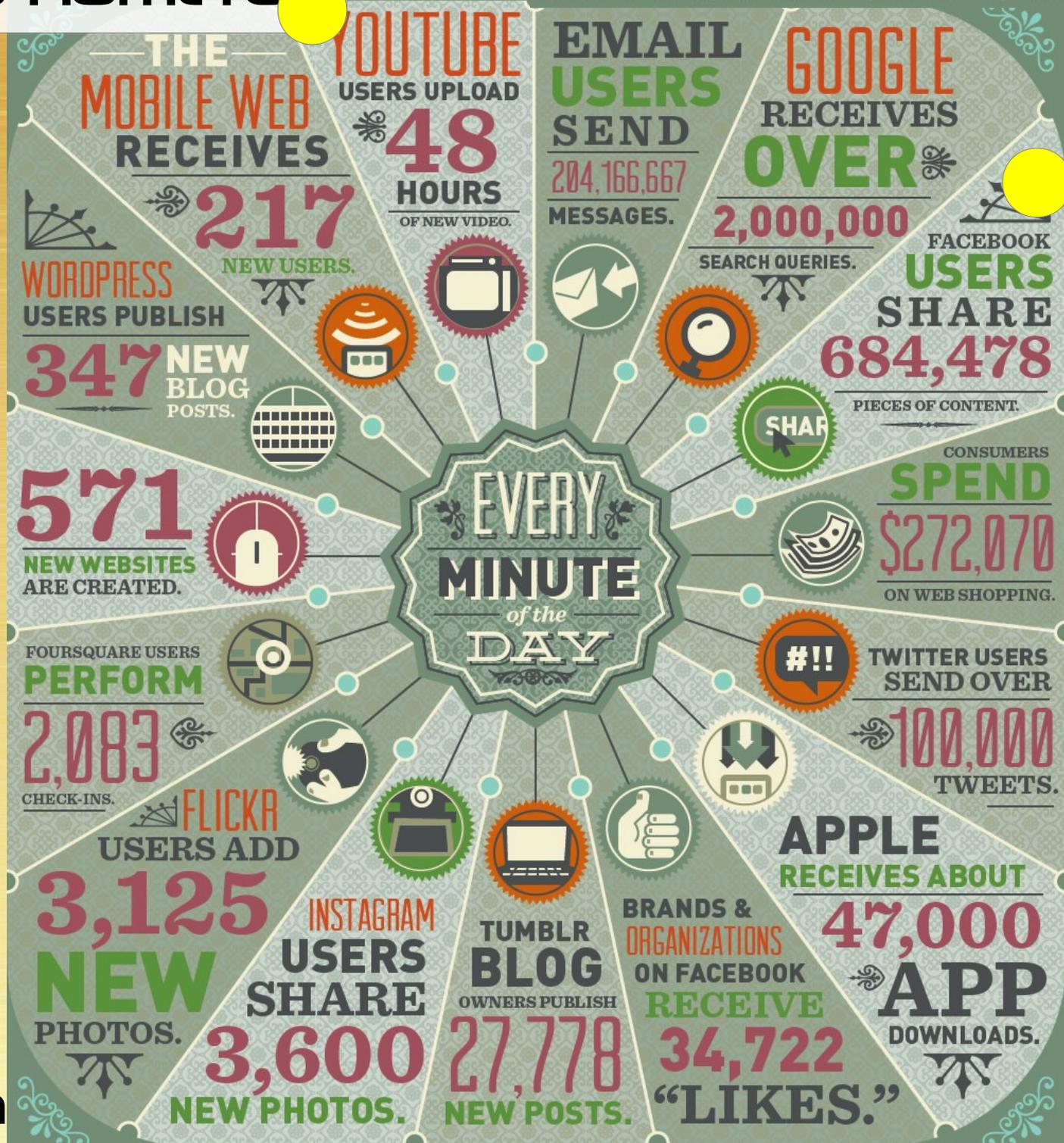
(120.000.000 datos útiles por minuto)

- 72M por hora
- 1.728M por día

0,6 Billones por año

2012

2.1 Billones de usuarios en Internet

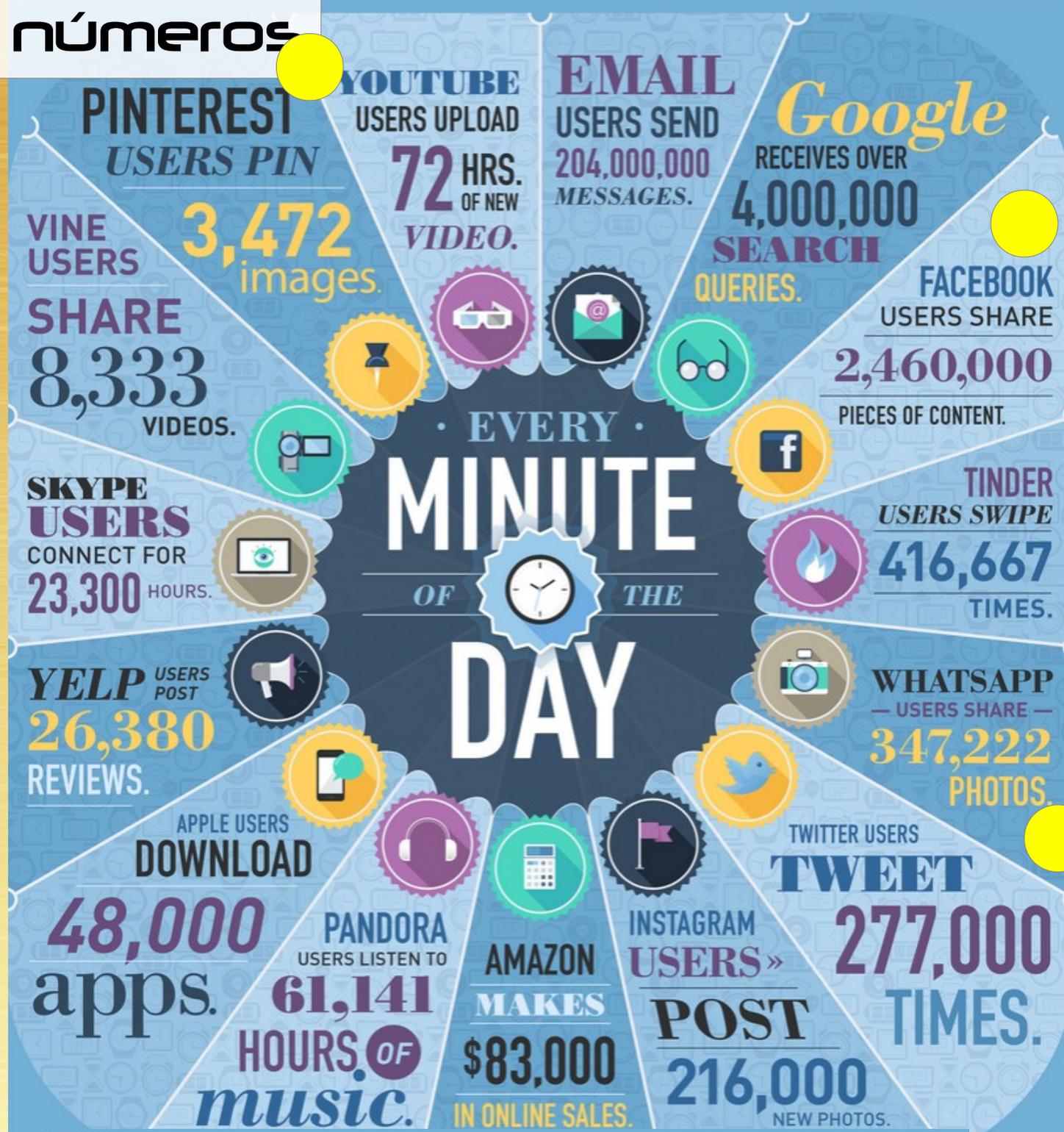


Grandes números

2014

De 2.1
a 2.4
billones

2012 - 2.1
2014 - 2.4



Juanmi

SOURCES: BITS.BLOGS.NYTIMES.COM, INTEL.COM, APPLE.COM, TIME.COM, DAILYMAIL.CO.UK, SKYPE.COM, STATISTICBRAIN.COM

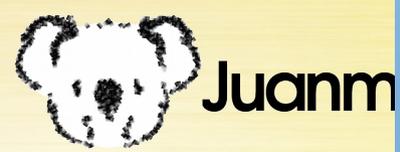
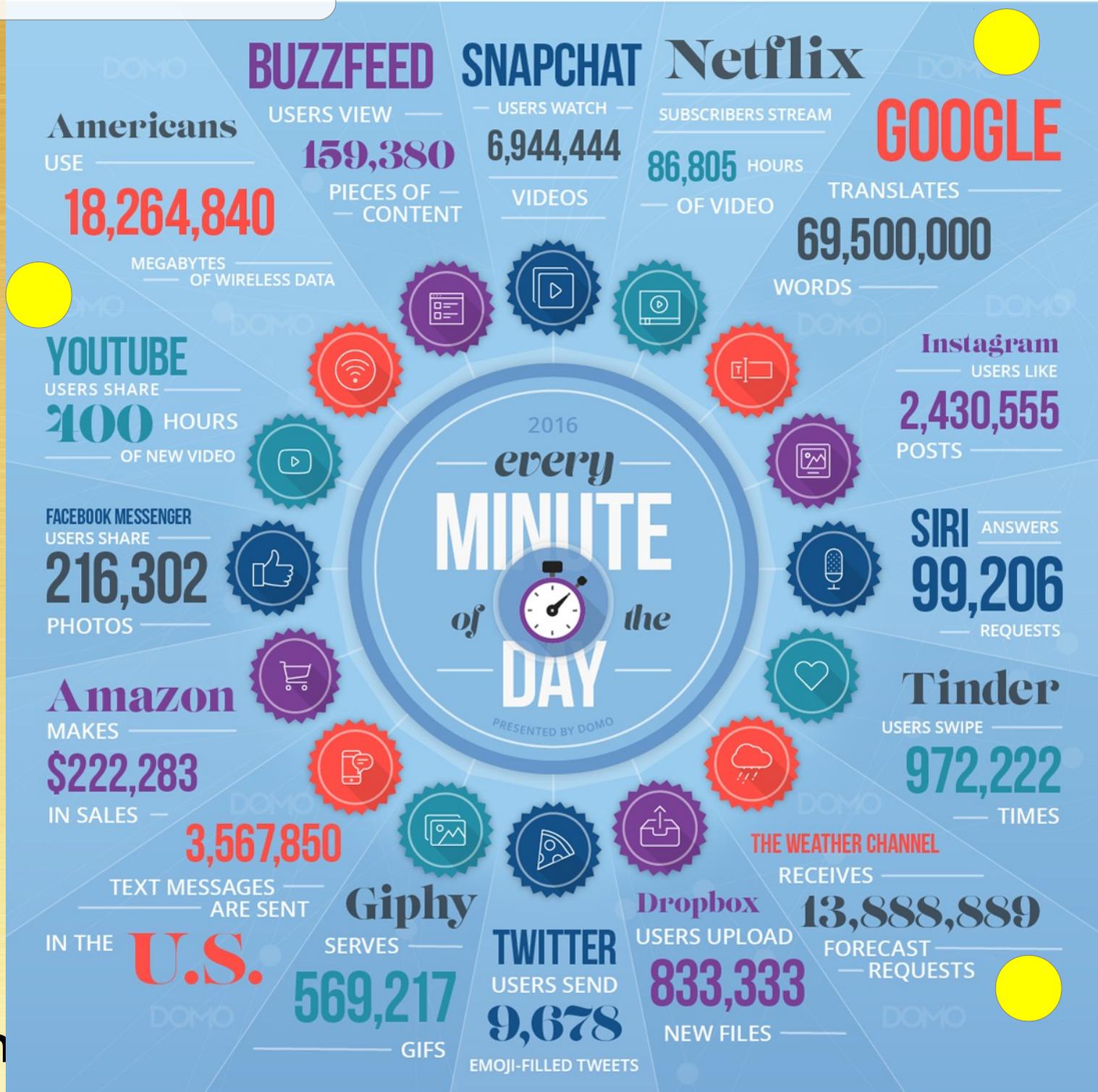
Data

Grandes números

2016

De 2.4
a 3.4
billones

2012 - 2.1
2014 - 2.4
2016 - 3.4

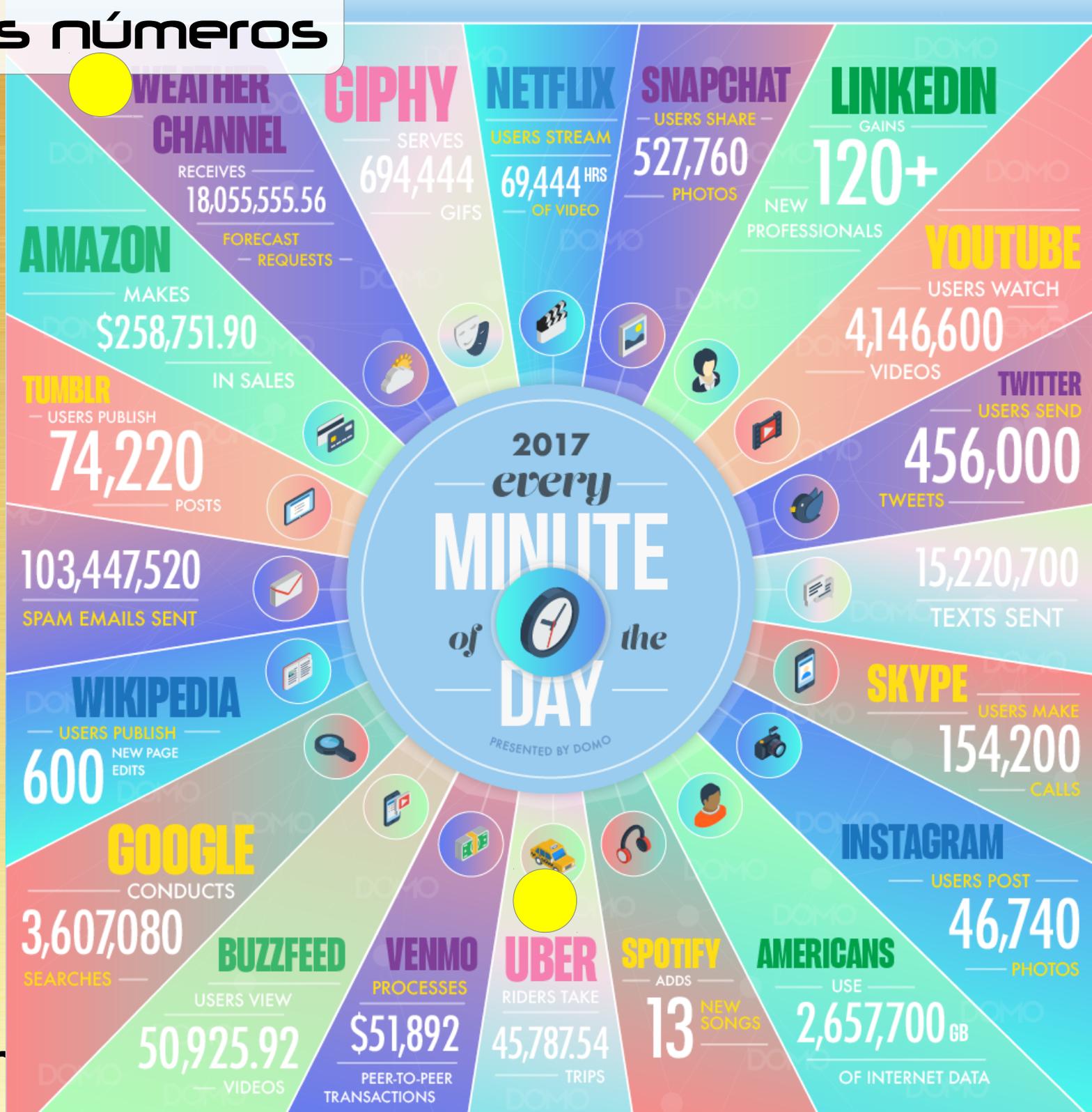


Grandes números

2017

De 3.4
a 3.7
billones

2012 - 2.1
2014 - 2.4
2016 - 3.4
2017 - 3.7

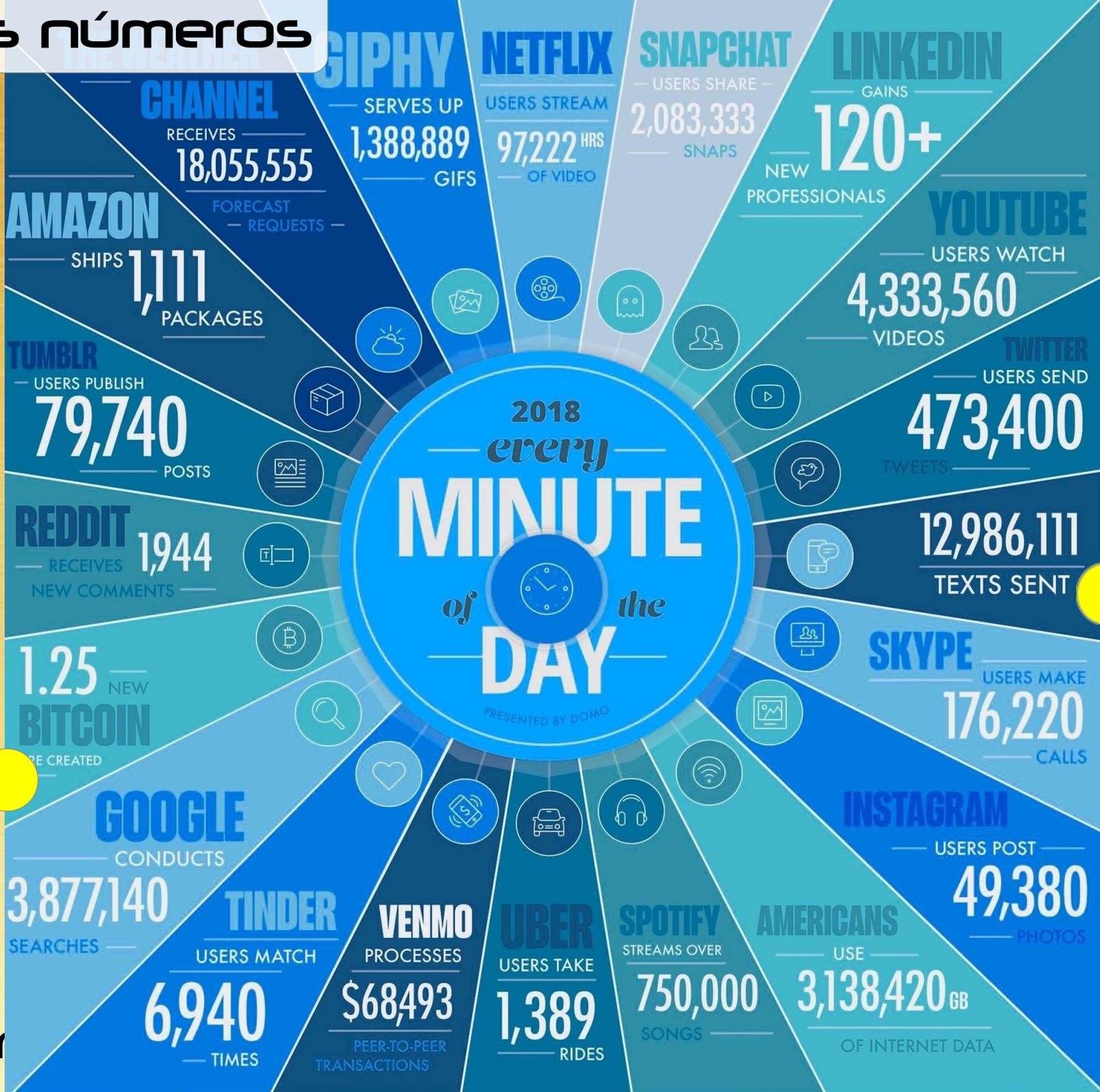


Grandes números

2018

De 3.4 a 3.7 billones

2012 - 2.1
2014 - 2.4
2016 - 3.4
2017 - 3.7
2017 - 4.4



ta

Usuarios de Internet:

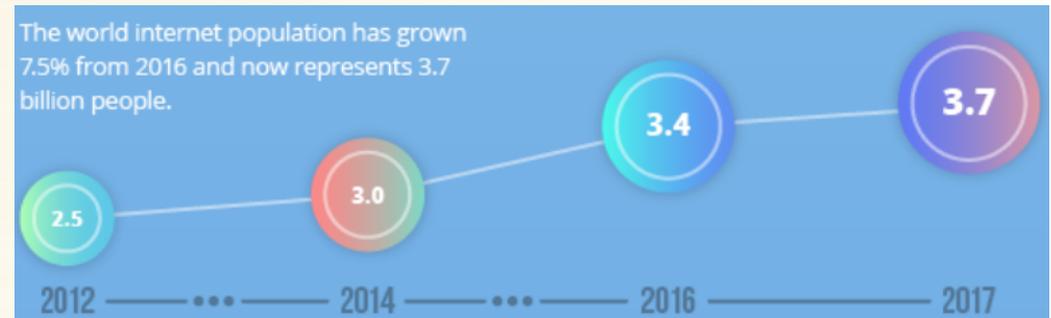
2012 – 2.1 Billones

2014 – 2.4 Billones

2016 – 3.4 Billones

2017 – 3.7 Billones

2018 – 4.4 Billones



Hoy: 86% del tráfico son videos

Toda la historia de películas juntas cruzan Internet **cada 3 minutos**

A lo largo de 1 año hablamos de 1 trillón de Gigabytes o ...

1 Zetabyte

IBM 2018: 2.5 Quintillones de bytes/día



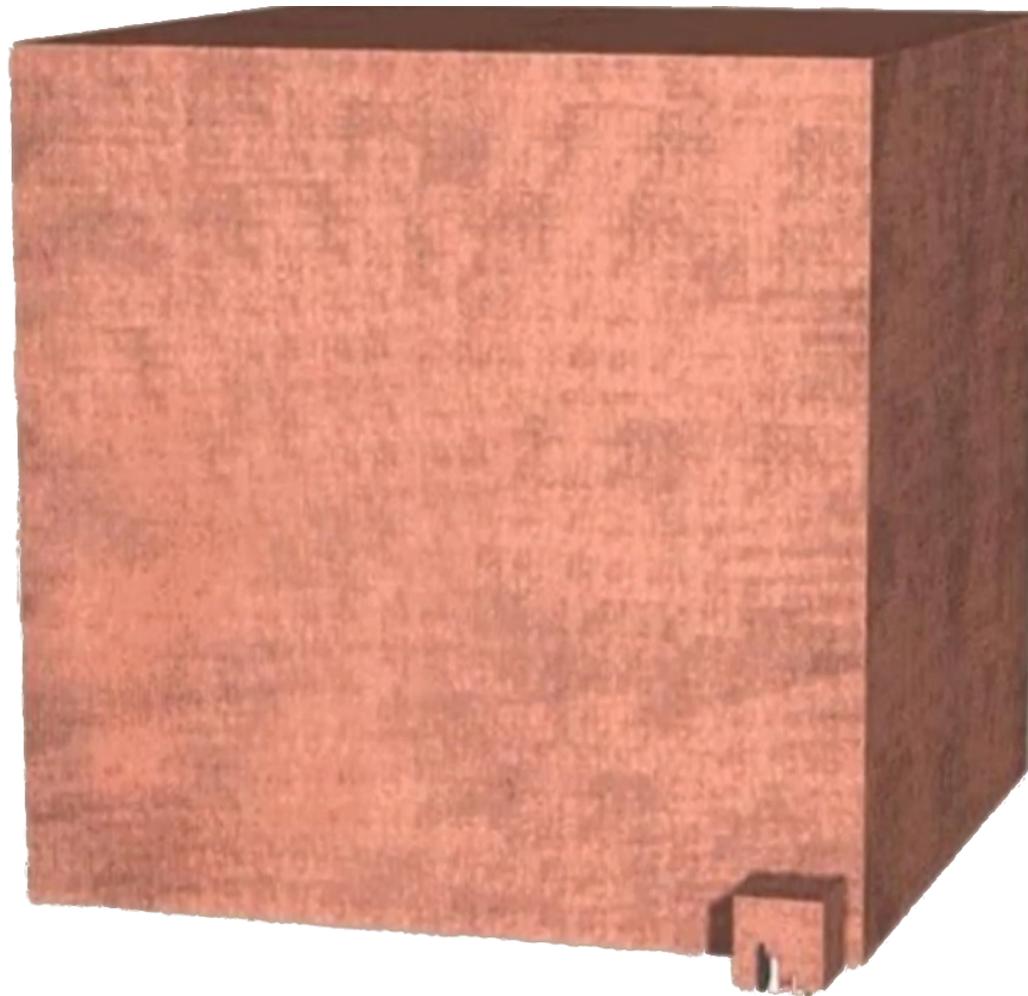
¿Cuanto es un quintillón?

- 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?



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



Eric Brewer (2000)

- C: Consistency** → **Consistencia**
- A: Availability** → **Disponibilidad**
- 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 que la mayoría **piensa** que tiene

Siempre se puede leer y escribir

A

C

P

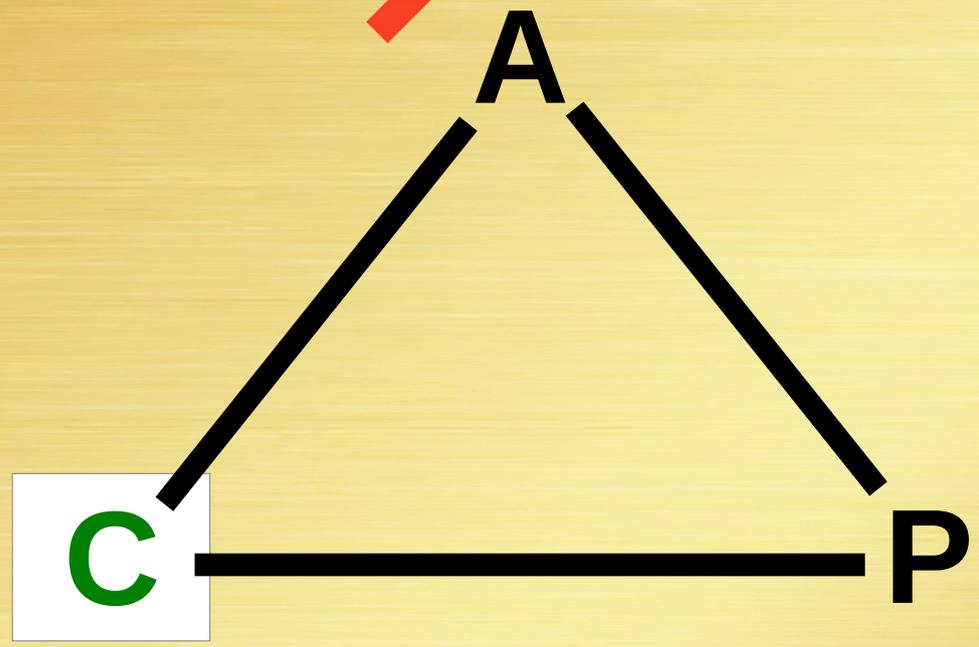
Todos los clientes ven siempre lo mismo

Los datos se pueden repartir en diferentes nodos

Teorema de CAP

Lo que la mayoría de **verdad** tiene

~~Siempre se puede leer y escribir~~



Todos los clientes ven siempre lo mismo

~~Los datos se pueden repartir en diferentes nodos~~

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: Estado flexible/maleable (Soft state)

E: Eventualidad del estado final

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

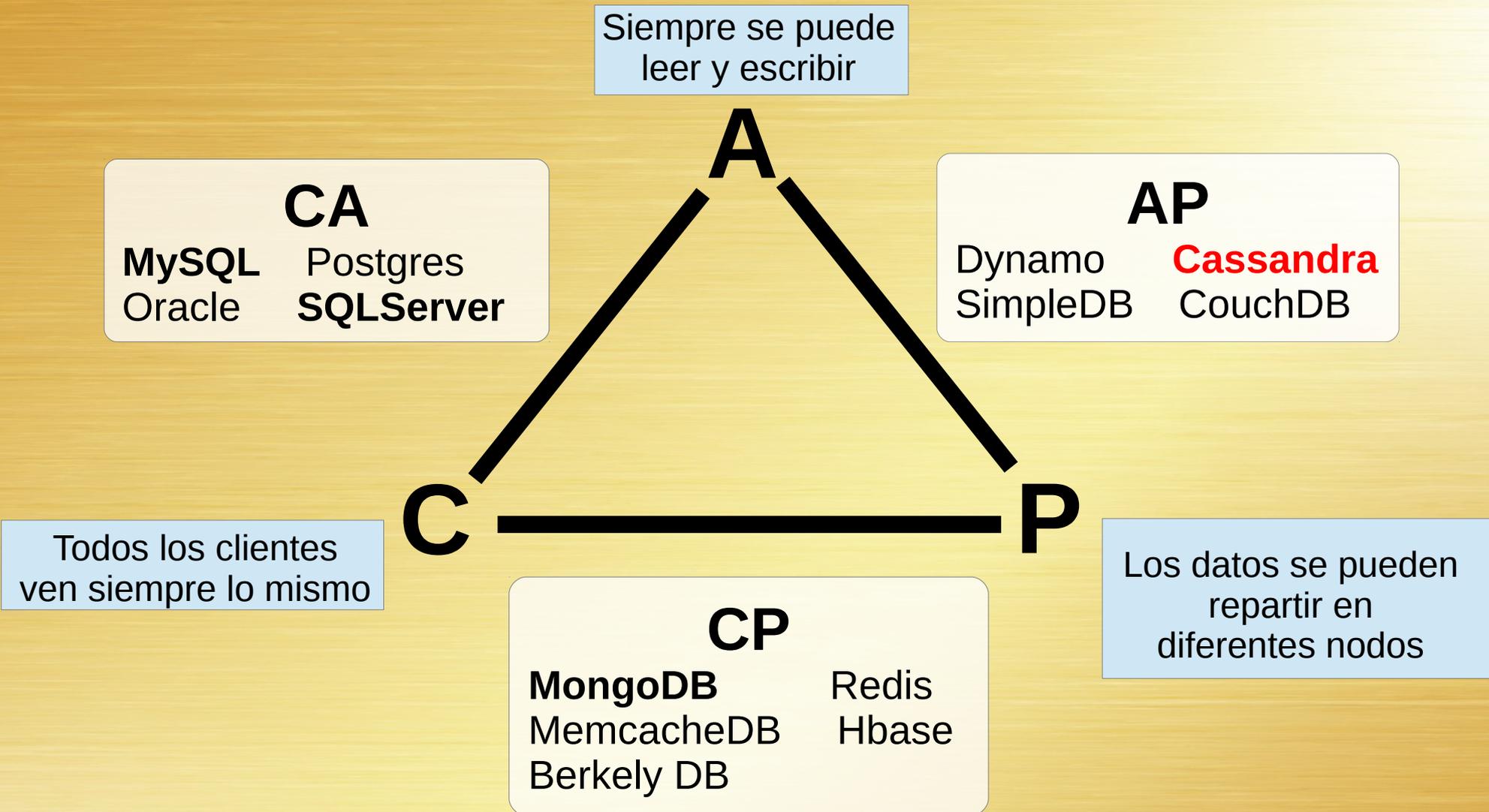
Métodos de resolución de conflictos

Evitan la entropía (control de versiones)

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

Garantizar la seguridad de una operación

Strong eventual consistency (SEC)



CA

MySQL
Postgres
Oracle
SQLServer

Aster Data
Greenplum
Vertica
Clickhouse

AP

Dynamo
Voldemort
Tokyo Cabinet
KAI

Cassandra
SimpleDB
CouchDB
Riak



CP

BigTable
Hypertable
HBase

MongoDB
Terrastore
Scalaris

Berkely DB
MemcacheDB
Redis

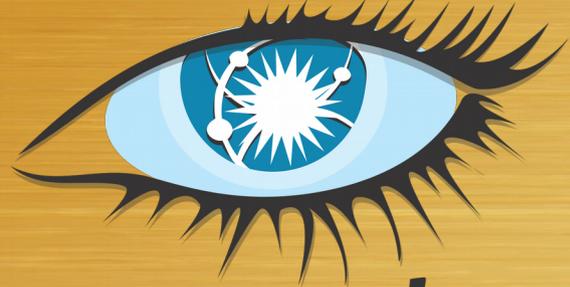




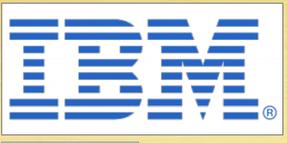
- Planificar antes de actuar, no te dejes llevar por las emociones. ¿Cuales son los objetivos?
- Diseño orientado a consulta: ¿Cuales son las preguntas a responder?
- ~~DRY~~: Vale repetir la información
- El almacenamiento no debe ser un coste significativo (HA, resiliencia, velocidad, esparcimiento,...)
- El motor de BD BigData va en f(proyecto, dato)

¿Quién usa Big Data?

Fuente: <http://planetcassandra.org/companies/>



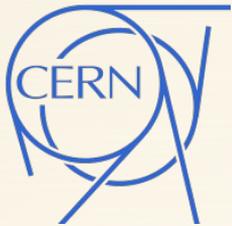
cassandra



Juanmi Taboada

BIG DATA

¿Quién usa Big Data?



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

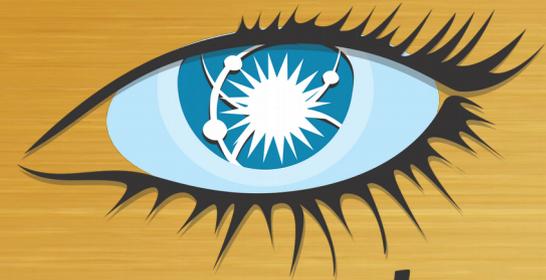
NETFLIX

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



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





cassandra

Disponibilidad continua

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

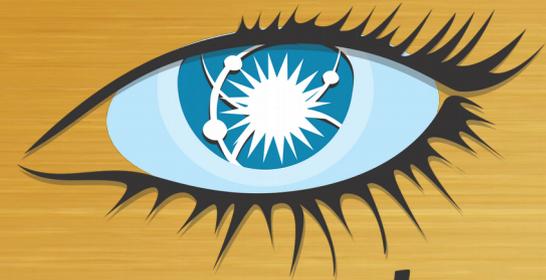
Escalabilidad lineal

- 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





cassandra

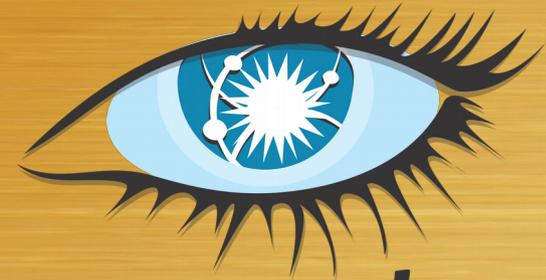
Gossip: mantiene la red informada

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



cassandra

Gossip:

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





cassandra

Particionador: encargado de “**esparcir**” los datos.

Desaconsejados:

Random: Hashes con MD5

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 distribución espacial.

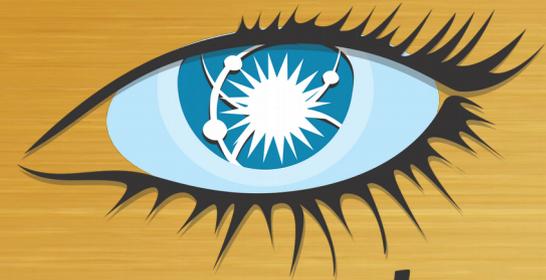




Planificar el deploy de un cluster

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



cassandra

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 (múltiples hits)
- Balanceadores de carga
- Falta de testing
- Bajo conocimiento de Linux





CQL

(Cassandra Query Language)

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;
```

CQL

```
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;
```

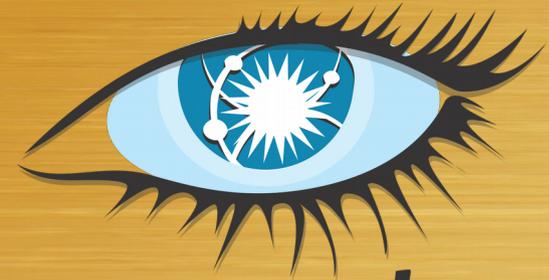
```
/* Counting Data */
```

```
SELECT COUNT(*) FROM myTable;
```

```
/* Deleting Data */
```

```
DELETE FROM myTable WHERE myField = 7;
```





cassandra

```
USE miBaseDatos;
```

```
/* Creando Tablas */
```

```
CREATE TABLE IF NOT EXISTS miTabla (id INT PRIMARY KEY);
```

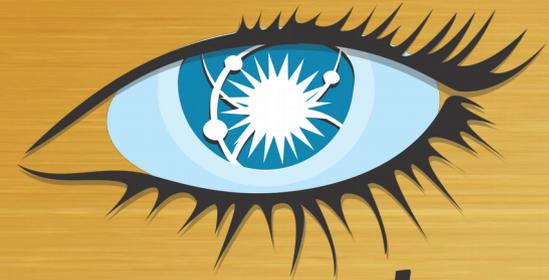
```
/* Cambiando Tablas */
```

```
ALTER TABLE miTabla ADD miCampo INT;
```

```
/* Creating Indexes */
```

```
CREATE INDEX miIndice ON miTabla (miCampo);
```





cassandra

/ 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;
```





cassandra

/ Crear un nuevo keyspace en CQL */*

```
CREATE KEYSPACE miBaseDatos WITH replication =  
{'class': 'SimpleStrategy', 'replication_factor': 1};
```

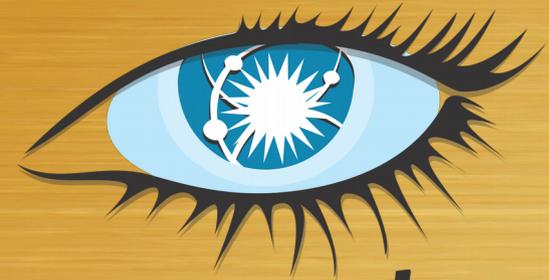
/ Crear una nueva base de datos en SQL */*

```
CREATE DATABASE miBaseDatos;
```

No existen:

JOIN, GROUP BY, y FOREIGN KEY

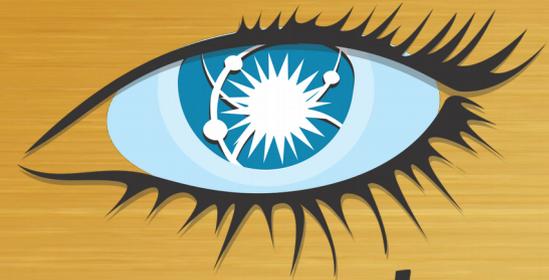




cassandra

- 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!





cassandra

```
/* 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**.





Prod1

Health: Normal 15, Medium 1, High 0, Down 0

Data Size: Total Size 321 GB, Avg. Node 20 GB, Standard Deviation +/- 502 MB

Alerts: 2

Ring View Explanation

Datacenter Analytics, Datacenter Cassandra, Datacenter Solr

machine2
101.202.203.106
Token: 576480752303423500
Status: Active
OS Load: 0.36
Data Size: 19.8 GB
Tracker: Active

Memory Usage

System: 7.6 GB, 8.4 GB

Heap: 3.2 GB, 861 MB

Storage Capacity

Load: 0.15

Status: Active

Token: 0

Gossip: Active

Thrift: Active

Native Transport: Active

Pending Tasks: 2

Running Tasks: No active tasks

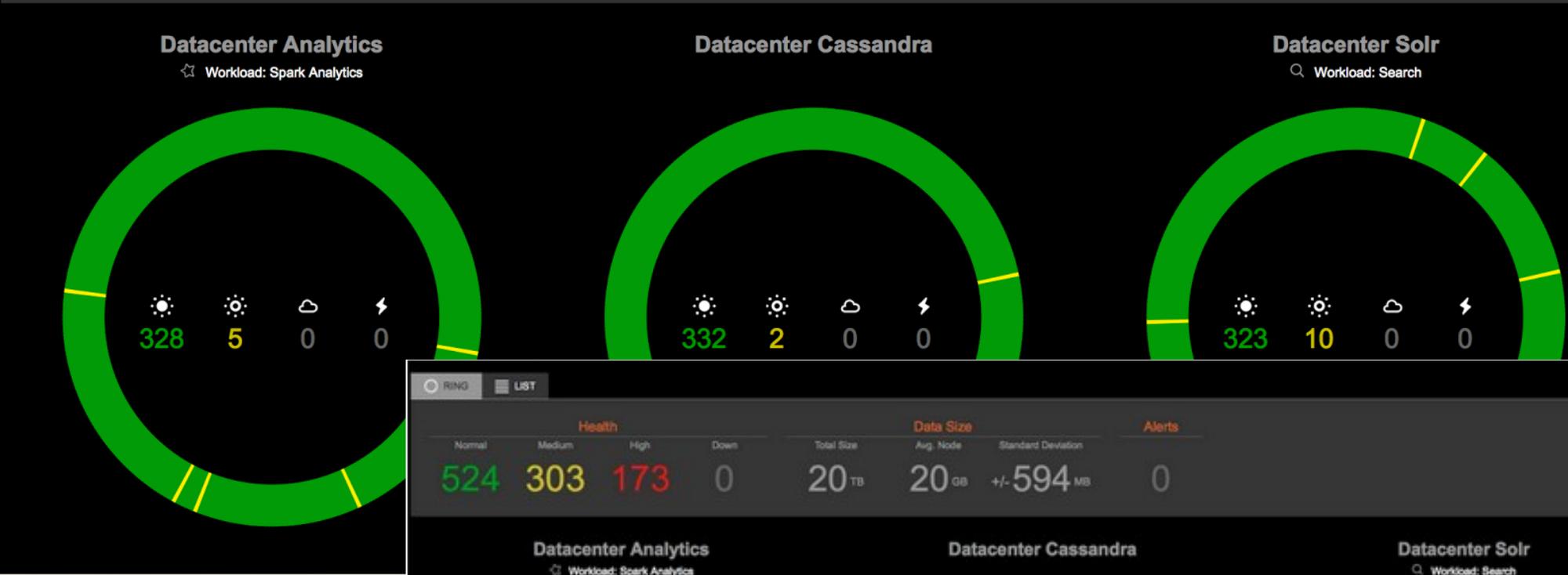
Streams: No active streams

Thread Pool Stats

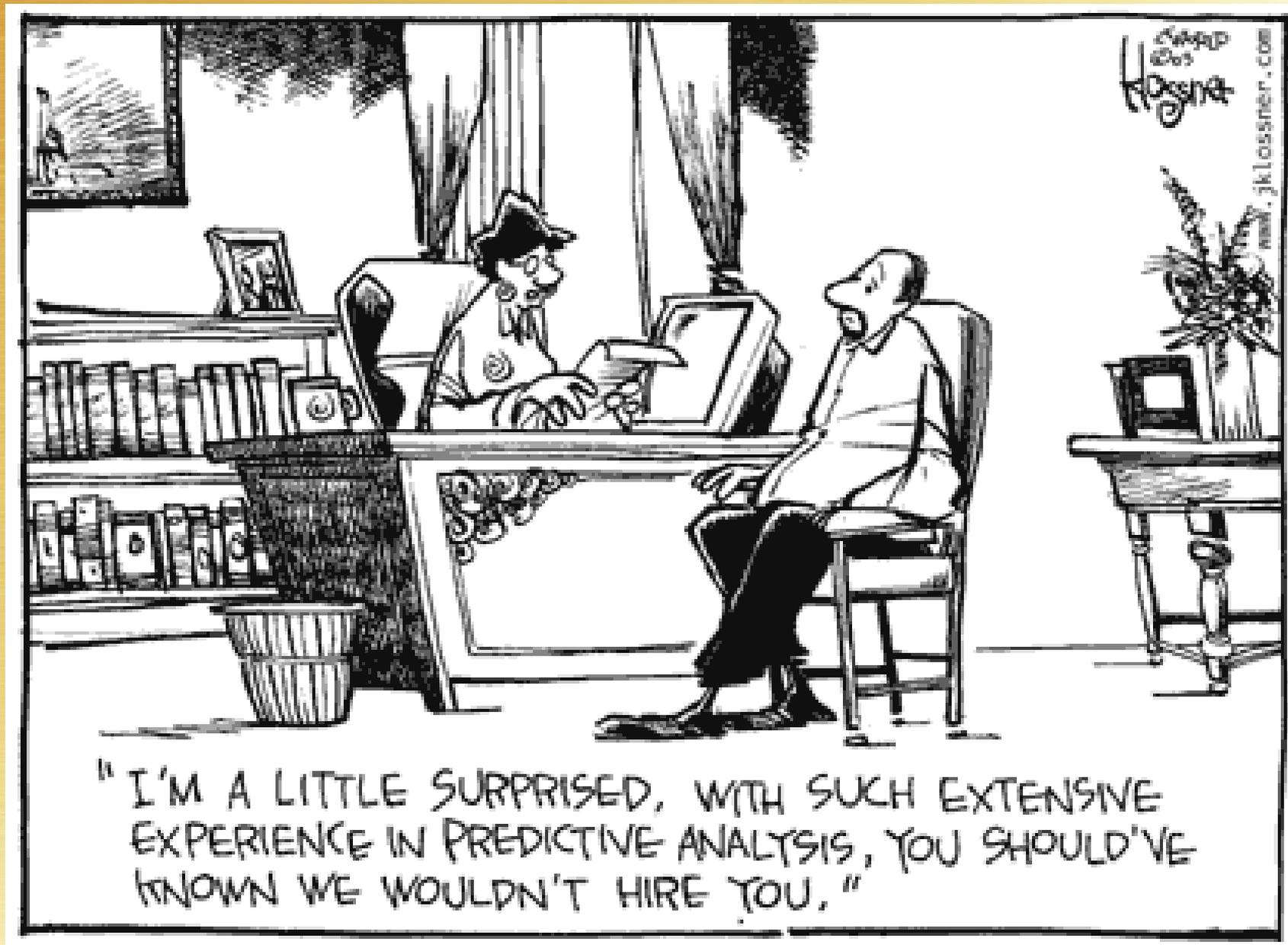
Name	Active	Pending	Completed	Blocked	Total Blocked
AntiEntropyStage	0	0	0	2	0
CommitLogArchiver	0	0	2	0	0



Health				Data Size			Alerts
Normal	Medium	High	Down	Total Size	Avg. Node	Standard Deviation	
983	17	0	0	20 TB	20 GB	+/- 593 MB	0



Una dosis de realidad

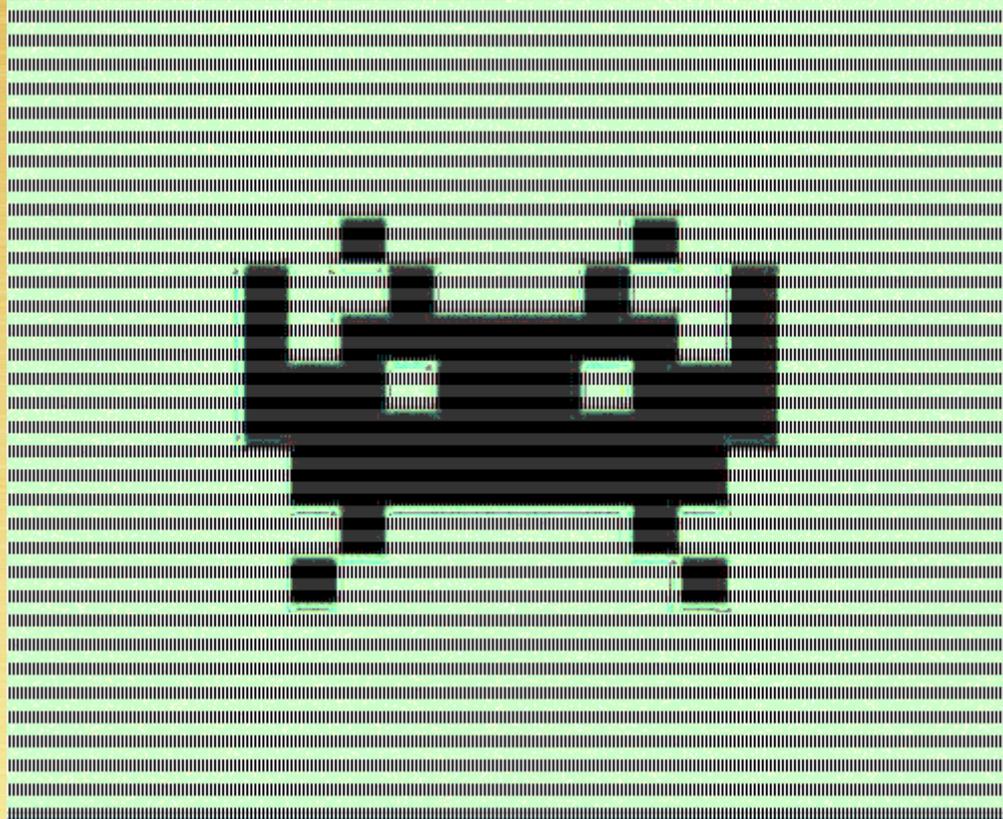


Una dosis de realidad

- Predicción de datos
- La tareas más ardua es filtrar y preparar los datos
- Desechar información para evitar males mayores
- PHP vs Python vs C, todo vale
- El rendimiento homogéneo de la BD es clave (red, Discos, CPU, RAM,..., tipo de partición)
- Hay bases de datos para diferentes tipos de datos (también existen bases de datos en memoria)
- Parar unos segundos son ciento de ventas perdidas



¿Dudas?



Muchas Gracias

Thank you - Dziękuję



Juanmi Taboada
@juanmitaboada
www.juanmitaboada.com

