

Sistemas Distribuidos

Diseño de Sistemas Distribuidos: Google

Rodrigo Santamaría

Diseño de Sistemas + Distribuidos

Introducción: Google como caso de estudio

Servicios

Plataforma

Middleware



Introducción

- Crear un sistema distribuido no es sencillo
 - **Objetivo:** obtener un sistema consistente que cumpla con los requisitos identificados
 - Escala, seguridad, disponibilidad, etc.
 - Elección del **modelo**
 - Tipo de fallos que asumimos
 - Tipo de arquitectura
 - Elección de la **infraestructura**
 - Middleware (RMI, REST, Kademlia, etc.)
 - Protocolos existentes (LDAP, HTTP, etc.)

+ Introducción

Google como caso de estudio

- Google es un ejemplo de un diseño distribuido exitoso
 - Ofrece búsquedas en Internet y aplicaciones web
 - Obtiene beneficios principalmente de la publicidad asociada
- **Objetivo:** *“organizar la información mundial y hacerla universalmente accesible y útil”*
- Nació a partir de un proyecto de investigación en la Universidad de Stanford, convirtiéndose en compañía con sede en California en 1998.
- Una parte de su éxito radica en el algoritmo de posicionamiento utilizado por su motor de búsqueda
 - El resto, en un sistema distribuido eficiente y altamente escalable

+ Introducción

Desafíos en Google

- Google se centra en cuatro desafíos
 - **Escalabilidad:** un sistema distribuido con varios subsistemas, dando servicio a millones de usuarios
 - **Fiabilidad:** el sistema debe funcionar en todo momento
 - **Rendimiento:** cuanto más rápida sea la búsqueda, más búsquedas podrá hacer el usuario -> mayor exposición a la publicidad
 - **Transparencia:** en cuanto a la capacidad de reutilizar la infraestructura disponible, tanto internamente como externamente (plataforma como servicio)

+ Introducción

Desafíos: escalabilidad

- Google ve estos desafíos en términos de tres dimensiones
 - *Datos*: el tamaño de Internet sigue creciendo (p.ej. por la digitalización de bibliotecas o los contenidos en redes sociales)
 - *Peticiones*: el número de usuarios sigue creciendo
 - *Resultados*: la calidad y cantidad de resultados sigue mejorando
- Logros
 - 1998: sistema de producción inicial
 - 2010: $88 \cdot 10^9$ millones de búsquedas al mes
 - Sin perder eficiencia: tiempo de consulta $< 0.2s$
- Ejemplo [Dean 2006]
 - Asumimos que la red está compuesta de $\sim 20 \cdot 10^9$ páginas
 - Cada una de 20KB \rightarrow tamaño total de 400TB
 - Un ordenador que lea 30MB/s tardaría 4 meses en explorarla
 - 1000 ordenadores lo harían en menos de 3h

+ Introducción

Desafíos: fiabilidad

- Google ofrece un nivel de fiabilidad del 99.9% en sus contratos de pago
 - Generalmente lo ha cumplido, pero en 2009, sufrió una caída de 100 minutos en Gmail durante un mantenimiento rutinario
 - Algunas caídas en servicios gratuitos, sobre todo Gmail
 - <https://en.wikipedia.org/wiki/Gmail#Outages>
- Para mantener la fiabilidad, es necesario anticipar los fallos (HW y SW) con una frecuencia razonable
 - Técnicas de detección de fallos
 - Estrategias para enmascarar o tolerar fallos -> principalmente mediante la replicación de la arquitectura física

+ Introducción

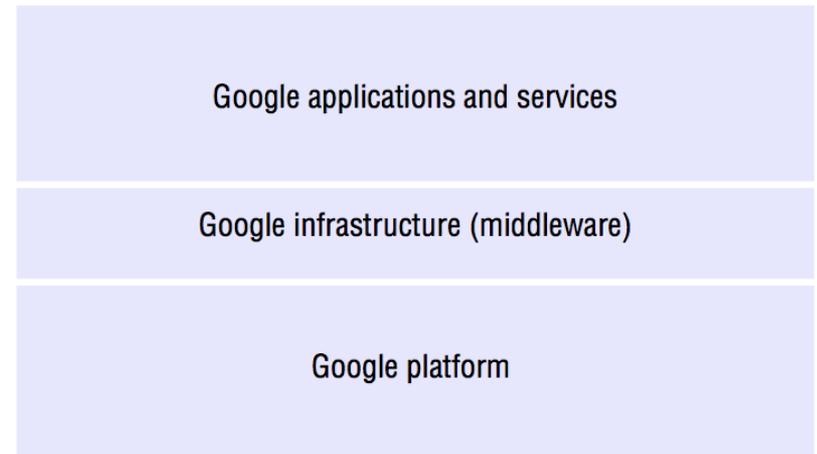
Desafíos: rendimiento

- Involucra a todas las fases de cada servicio
 - Por ejemplo, en cuanto al motor de búsqueda
 - Involucra a crawling, indexing y ranking
 - Objetivo: operaciones de búsqueda $< 0.2s$
- Involucra a todos los recursos subyacentes
 - Red
 - Almacenamiento
 - Procesamiento

+ Introducción

Desafíos: transparencia

- Requisito fundamental para el uso de la plataforma de Google como servicio, entendida como:
 - Extensibilidad
 - Soporte para el desarrollo de nuevas aplicaciones
- Para ello Google desarrolla su propia infraestructura (middleware) entre la arquitectura física y las aplicaciones y servicios



+ Google

Principios de diseño

- **Simplicidad:** 'cada API tiene que ser tan sencilla como sea posible y no más'
 - Aplicación de la navaja de Occam
- **Rendimiento:** 'cada milisegundo cuenta'
 - Medición de los costes de operaciones primitivas (acceso a disco y memoria, envío de paquetes, etc.)
- **Testeo:** 'si no se ha roto, es que no lo has probado lo suficiente'
 - Régimen estricto de pruebas y monitorización

Diseño de Sistemas + Distribuidos

Introducción

Servicios

- Motor de búsqueda
- Cloud computing

Plataforma

Middleware

+ Motor de búsqueda

- Objetivo:
 - Dada una consulta, obtener una lista ordenada de resultados relevantes, a partir de una búsqueda en la red.
- Subsistemas del motor:
 - *Crawling*: obtener información
 - *Indexing*: procesar información
 - *Ranking*: clasificar información

+ Motor de búsqueda

Crawling

- **Tarea:** localizar y recuperar contenidos de la Web y pasarlos al sistema de indizado.
- **Ejecución:** servicio software *Googlebot*
 1. Lee recursivamente una página web
 2. Recolecta todos sus enlaces
 3. Planifica posteriores operaciones de crawling en dichos enlaces
 - Esta técnica (*deep searching*) permite penetrar en prácticamente todas las páginas indexadas de la Web*
- La ejecución se realiza en batería. En el pasado, el Googlebot se ejecutaba una vez cada pocas semanas
 - No es suficiente para páginas de noticias o con otro tipo de contenido cambiante (blogs, redes sociales, etc.)

*Todas las páginas indexadas de la web **no** son toda la web
https://en.wikipedia.org/wiki/Deep_web

+ Motor de búsqueda

Crawling

- Con la arquitectura de búsqueda *Caffeine*, se introduce en 2010 una nueva filosofía de crawling/indexing
 - El crawler está en funcionamiento *continuo*, y tan pronto como se explora una página, se realiza su indexado.
- Este modo reduce la antigüedad de los índices en un 50%
 - Se pasa de un procesamiento por lotes global a un procesamiento continuo **incremental**.
 - Es una estrategia clásica de mejora (compresión de archivos, copias de seguridad, sistemas de versiones ...)

+ Motor de búsqueda

Indexing

- **Tarea:** producir un índice de contenidos de la red similar al de un libro (título, autor, tema, etc.)
- Ejecución: *indizado* inverso de palabras
 - Dada una página web (varios formatos: html, pdf, doc), se identifican características textuales clave: posición, tamaño de letra, capitalización.
 - También se realiza un índice de enlaces encontrados en la página
- Utilizando este índice, se reduce el número de páginas candidatas de miles de millones a unas decenas de miles, según el poder discriminativo de las palabras buscadas.

+ Motor de búsqueda

Ranking

- **Tarea:** clasificar los índices en función de su relevancia
- **Ejecución:** algoritmo PageRank
 - Basado en los sistemas de ranking de publicaciones científicas: un artículo científico es importante si ha sido citado por otros colegas del área.
 - Análogamente, una página es importante si ha sido enlazada por un gran número de páginas.
 - También tiene en cuenta factores relacionados con la proximidad de la búsqueda a las palabras claves de la página obtenidas en el indizado inverso

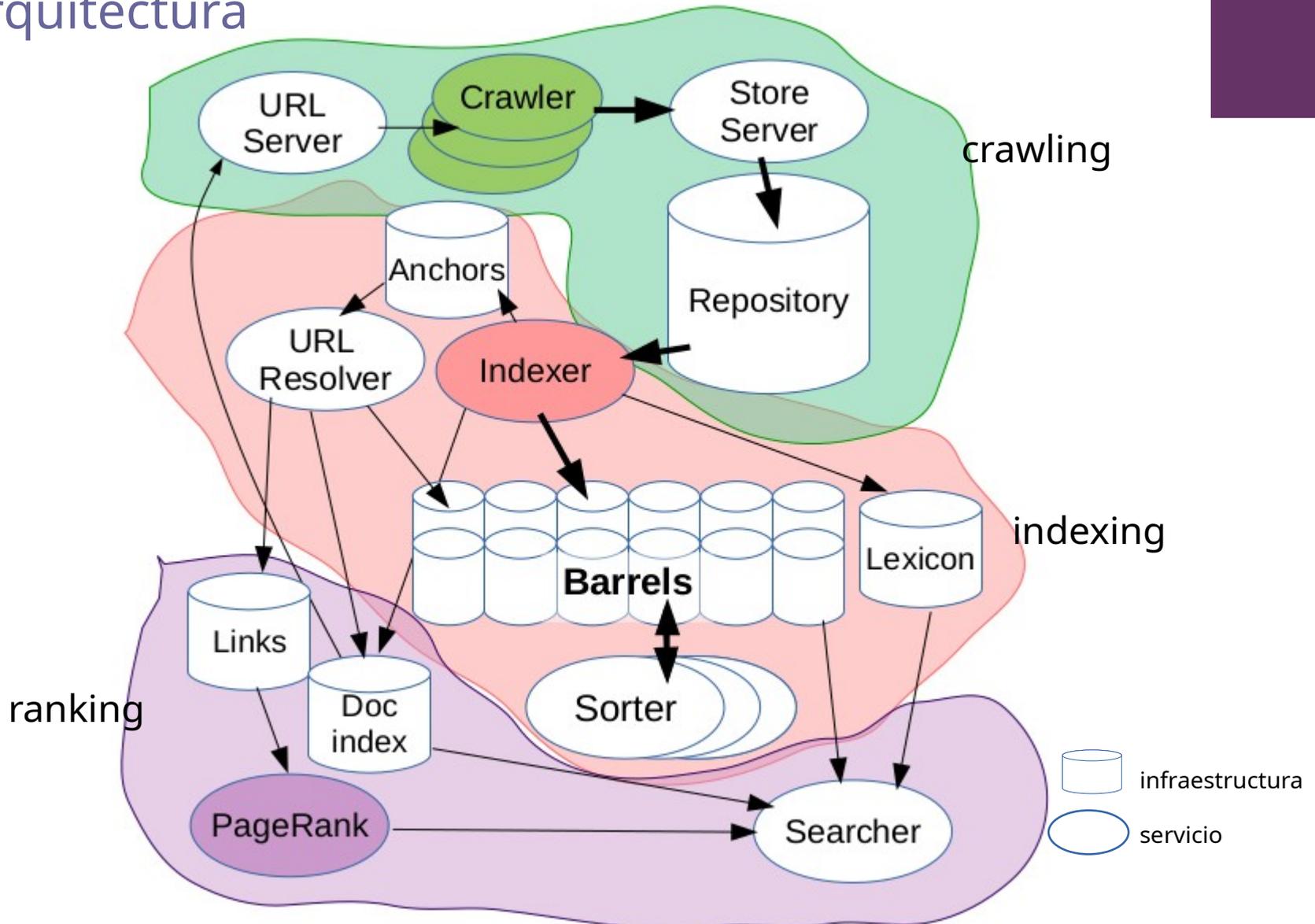
+ Motor de búsqueda

Ejemplo de búsqueda

- Imaginemos que buscamos “distributed systems book”
 - El buscador seleccionará páginas que contengan las tres palabras en sus índices
 - Preferentemente en el título o alguna sección importante (p. ej. listas de libros de sistemas distribuidos)
 - El buscador priorizará aquellas con un alto número de enlaces, poniendo primero las más enlazadas, y desde páginas más ‘importantes’:
 - *amazon, wikipedia, bookdepository* son priorizadas por su relevancia en la red.

+ Motor de búsqueda

Arquitectura



+ Motor de búsqueda

Arquitectura

- *Repository*: datos web recogidos por el crawler, comprimidos
- *Barrel(s)*: almacén de *hits*, es decir, entradas que contienen el identificador del documento, la palabra encontrada en el documento, su tamaño de letra y capitalización
- *Sorter*: reordena el barrel por identificadores de palabra para generar el índice invertido
- *Anchors*: información sobre los enlaces en los documentos
- *URL resolver*
 - Preparación de los enlaces (*links*)
 - Obtención de nuevos enlaces para alimentar al Crawler (*Doc index*)

+ Motor de búsqueda

Arquitectura

- Los detalles específicos de esta arquitectura han variado, pero los servicios clave (crawling, indexing, ranking) siguen siendo los mismos.
- Modificaciones en
 - Almacenamiento optimizado
 - Bloques de comunicación más reutilizables
 - Procesamiento por lotes a continuo

+ Motor de búsqueda

Críticas

- Peligro de autocomplacencia [Maurer, 2007]
 - Google nos dice qué es importante -> no miramos más allá
- ¿importancia = n° enlaces?
 - Aquello más enlazado es más famoso, no necesariamente más relevante o interesante para el usuario
 - Discriminación de los sitios nuevos
- Peligro de manipulación
 - Neutralidad: Google prioriza sus servicios en las búsquedas [NexTag, 2011]
 - Google bombing: enlazar con un determinado texto a una página
 - Enlaces 'miserable failure' a páginas relacionadas con George Bush
 - Search Engine Optimization (SEO): modificar los contenidos de una página en base a los buscadores y los usuarios (~marketing)

+ Motor de búsqueda

Un cambio de prioridades

- En 2020, Benjamin Gomes dejó de ser el director del servicio de búsqueda tras 20 años, siendo reemplazado por Prabhakar Raghavan, especialista en marketing*
 - En un cruce de correos, Gomes mostró su preocupación: “*we are getting too close to money*” ([enlace](#))
- El cambio fue exitoso en económicamente, pero degradó sensiblemente el producto:
 - Eliminación de comprobaciones de páginas poco fiables
 - En particular, retirada del ‘Penguin’ update
 - Camuflaje de las páginas promocionadas

* <https://www.wheresyoured.at/the-men-who-killed-google/>

+ Google

Otros servicios: Cloud computing

- Modificación de un cliente ligero que permite usar archivos y aplicaciones remotas sin perder autonomía.
- Software como servicio: Google Apps
 - Alternativa a las suites de escritorio:
 - Gmail, Google Drive, Google Sites, Google Calendar
 - Google Maps, Google Earth
- Plataforma como servicio: Google App Engine
 - Pone a disposición del usuario parte de la infraestructura distribuida de Google para sus Google Apps y su servidor de búsqueda
 - Para que pueda desarrollar sus propias aplicaciones web mediante su plataforma
 - Mediante el uso de APIs

Diseño de Sistemas + Distribuidos

Introducción

Servicios

Plataforma

- Modelo físico
- Componentes

Middleware

+ Arquitectura

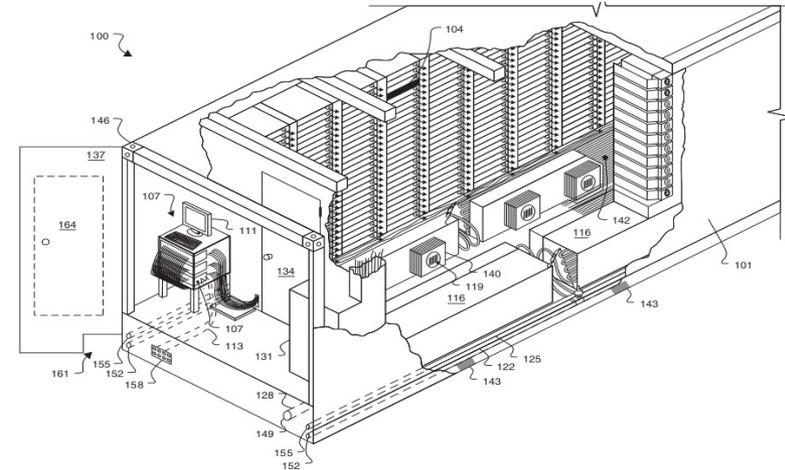
Modelo físico

- Utilización de gran número de PCs como base para construir un entorno para computación y almacenamiento distribuido
 - Cada unidad de PC utilizada cuesta entorno a 1000\$
 - Típicamente tiene 2TB de disco y 16GB de DRAM
 - Corre con una versión reducida del kernel de Linux
- Modelo de fallos
 - Al utilizar PCs 'de coste', hay un riesgo de fallos
 - Según [Hennessy and Patterson, 2006. Figura 6.1]
 - ~14 máquinas tienen fallos software cada día (reinicios)
 - 2-3% de los PCs tienen un fallo hardware al año (el 95% son fallos en los discos o en la DRAM) -> Un fallo HW por cada 10 fallos SW.
 - Se diseñarán estrategias de tolerancia de fallos

+ Arquitectura

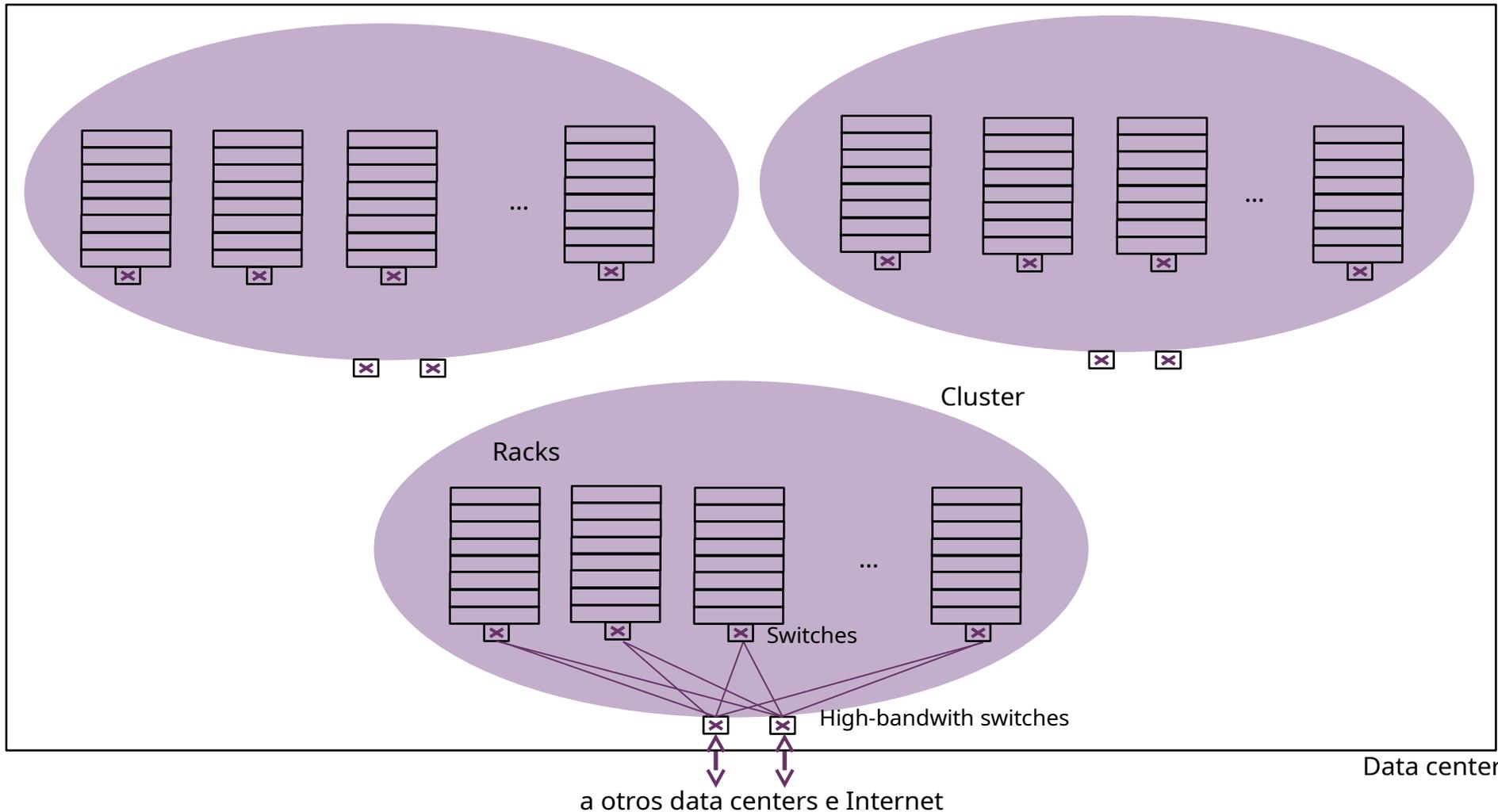
Componentes [Hennesy and Patterson, 2006]

- Rack
 - Entre 40 y 80 PCs
 - Switch Ethernet que provee conexión en el rack y hacia fuera
- Cluster
 - 30+ racks
 - 2 switches de banda ancha (redundancia)
 - Unidad básica de gestión
- Data Center
 - Decenas repartidos por el mundo
 - Cada uno hospeda uno o más clusters.



+ Arquitectura

Componentes



+ Arquitectura

Data Centers

- 12 dedicados, 24 compartidos (2008)



<http://royal.pingdom.com/2008/04/11/map-of-all-google-data-center-locations/>

<http://www.google.com/about/datacenters/locations/index.html>



Arquitectura

Capacidad de almacenamiento

- Un PC -> 2 terabytes
- Un rack de 80PCs -> 160 terabytes
- Un cluster de 30 racks -> 4.8 petabytes
- Asumiendo unos 200 clusters entre todos los data centers -> 960 petabytes ~ 1 exabyte (10^{18} bytes)

Diseño de Sistemas + Distribuidos

Introducción

Servicios

Plataforma

Middleware

- Comunicación
- Almacenamiento
- Computación

+ Middleware

- Para implementar los servicios ofertados a partir de la plataforma distribuida disponible, hace falta toda una infraestructura (middleware) intermedia
- Por lo general, Google implementa sus propias soluciones middleware en vez de utilizar estándares
 - Son muy parecidas a las soluciones existentes
 - Pero están optimizadas para los requisitos de servicio y las características de la plataforma
- En esta sección veremos por encima cómo es este middleware en términos de comunicación, almacenamiento y computación distribuida

+ Middleware

Secciones

- **Paradigmas de comunicación:** invocación remota y multidifusión
 - *Protocol buffers:* ofrecen un formato de serialización y comunicación común
 - *Publish-subscribe:* servicio para la disseminación de eventos
- **Datos y coordinación:** almacenamiento y acceso coordinado a datos
 - *GFS:* sistema de archivos distribuido para grandes volúmenes de datos
 - *Chubby:* almacenamiento de volúmenes pequeños de datos y coordinación
 - *Bigtable:* base de datos distribuida construida sobre GFS/Chubby
- **Computación distribuida:** paralelización sobre la arquitectura física
 - *MapReduce:* computación distribuida sobre conjuntos grandes de datos
 - *Go/Sawzall:* lenguaje/biblioteca para computación distribuida

+ Middleware

Comunicación – protocol buffers

- Siguiendo el principio de simplicidad, Google adopta un servicio de invocación remota mínimo y eficiente
- Para ello utiliza los búferes de protocolo (*protocol buffers*), un mecanismo general para
 - Almacenamiento/marshalling, más sencillo que XML
 - Comunicación, intercambiándose mediante RPC
- Los búferes de protocolo son neutrales respecto a
 - Plataforma
 - Lenguaje
 - Protocolo RPC

+ Middleware

Comunicación – protocol buffers

- Utilizan un lenguaje sencillo de especificación de *mensajes*
- Cada *mensaje* se parece a un objeto:
 - Con campos enumerados
 - Secuencialmente
 - Y etiquetados
 - Política de existencia (repetido, requerido, opcional)
 - Con mensajes anidados
- Para cada campo, el compilador genera métodos de gestión
 - exists, clear, set, get

```
message Book {  
    required string title = 1;  
    repeated string author = 2;  
    enum Status {  
        IN_PRESS = 0;  
        PUBLISHED = 1;  
        OUT_OF_PRINT = 2;  
    }  
    message BookStats {  
        required int32 sales = 1;  
        optional int32 citations = 2;  
        optional Status bookstatus = 3 [default = PUBLISHED];  
    }  
    optional BookStats statistics = 3;  
    repeated string keyword = 4;  
}
```

+ Middleware

Comunicación – protocol buffers

- Los búferes de protocolo son más ligeros que XML
 - No son tan abundantes en tags
 - 30-10% del tamaño en XML
 - Métodos de acceso y numeración secuencial
 - 10 a 100 veces más rápidos en operación y acceso
- La comparación no es del todo justa
 - XML tiene un propósito general
 - Tiene una semántica mucho más rica
 - Busca la interoperabilidad entre todo tipo de sistemas



Middleware

Comunicación – protocol buffers

- Los búferes de protocolo se intercambian mediante RPC con una sintaxis adicional para generar interfaces y métodos

```
service SearchService {  
    rpc Search (RequestType) returns (ResponseType);  
}
```

- Interfaz de un servicio de búsqueda *SearchService*
 - Con una operación remota *Search* que toma un parámetro del tipo *RequestType* y retorna uno de tipo *ResponseType*
- Sólo se puede enviar un parámetro y retornar otro
 - Simplicidad
 - Filosofía tipo REST: operaciones simples, manipulación de recursos
- El compilador se encarga de generar los stubs para la comunicación RPC a partir del código

+ Middleware

Comunicación – publish-suscribe

- Ampliación de los protocol buffers para tareas de diseminación de eventos
 - RPC tiene un rendimiento bajo para este tipo de tarea
 - Necesidad de conocer la identidad de todos los destinatarios
- Google no ha hecho públicos sus detalles pero básicamente es un sistema de notificación de eventos
 - Un servicio publica los eventos que genera
 - Los suscriptores expresan su interés en determinados eventos
 - El sistema asegura la entrega de las notificaciones de eventos del servicio a sus suscriptores

+ Middleware

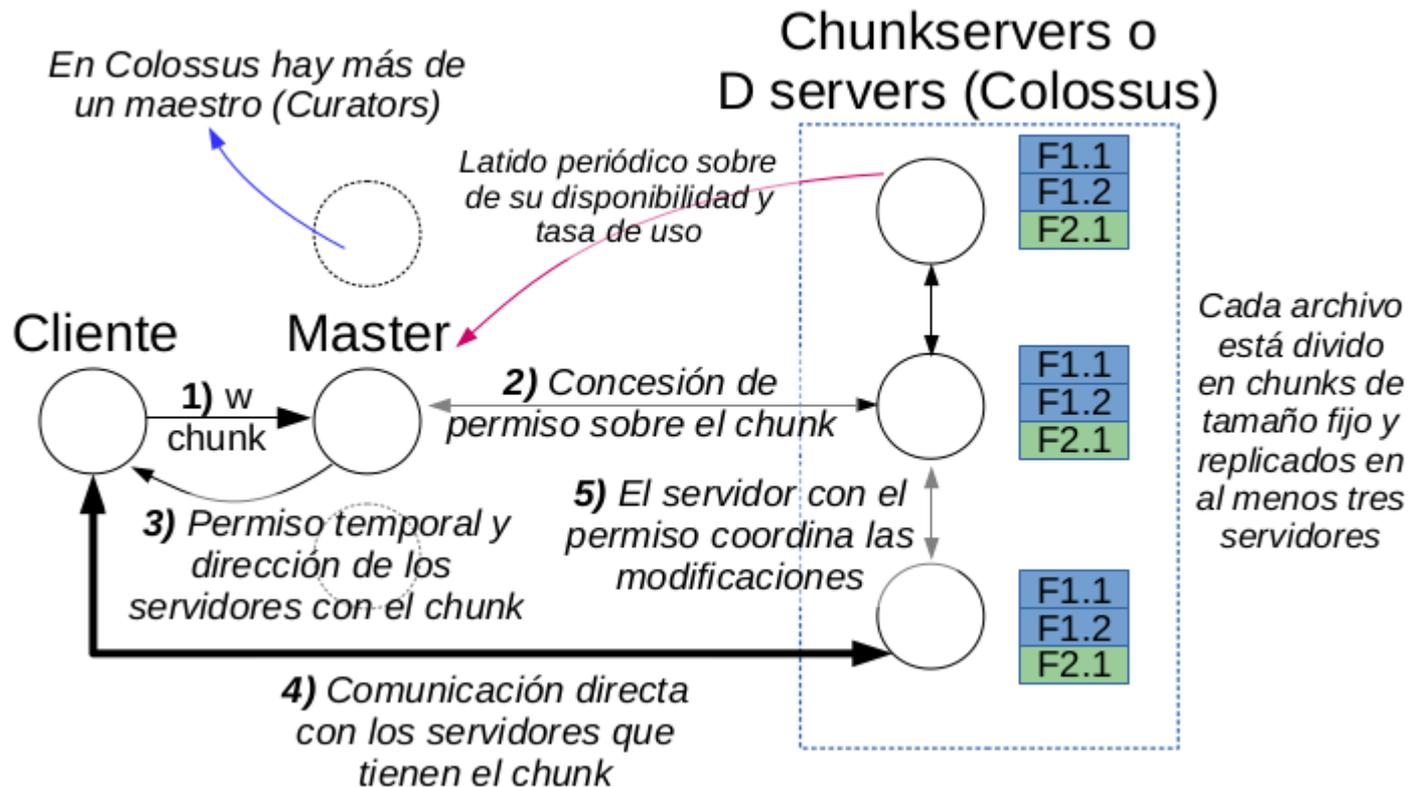
Almacenamiento - GFS

- Google File System es un sistema de archivos distribuido
 - Similar a otros de propósito general como NFS o AFS
 - Pero diseñado para los requisitos de Google:
 - Arquitectura física basada en clusters → replicación
 - Archivos de gran tamaño → chunks
 - Consultas secuenciales → permisos
 - Acceso concurrente alto → secciones críticas
 - Cada archivo se divide en trozos (**chunks**) que se almacenan **replicados** al menos tres veces (más según nivel de uso)
 - Su uso se regula con un sistema de **permisos**, similar a secciones críticas actualizadas mediante **latidos** periódicos

+ Middleware

Almacenamiento – GFS/Colossus

- Colossus sustituye a GFS en 2010, replicando nodos maestros y usando chunks y estructuras de datos más pequeñas (e.g. Bigtable).



+ Middleware

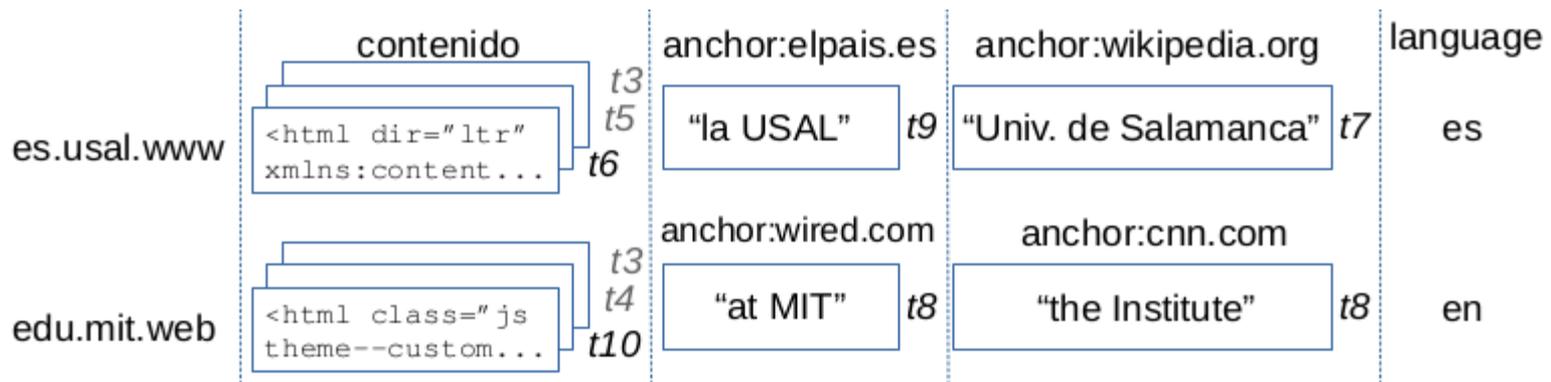
Almacenamiento – Bigtable

- GFS es un sistema de almacenamiento de archivos grandes
 - Pero 'planos', no indexados, accedidos secuencialmente
- Necesidad de acceso estructurado a los datos
 - Pero una base de datos distribuida es difícil de escalar y tiene bajo rendimiento, y no se necesita toda su funcionalidad
- Bigtable es un tipo de BD no convencional, donde las entradas de la tabla pueden tener un número y significado distinto de columnas.
 - Cada tabla suele estar en el rango de los terabytes
 - Las tablas estarán almacenadas mediante GFS/Colossus

+ Middleware

Almacenamiento - Bigtable: interfaz

- Cada tabla se accede a través de tres dimensiones
 - Fila: indica el objeto
 - P. ej. para el indizado inverso será la página web invertida para optimizar las búsquedas
 - Columna: indica el atributo - variable para cada entrada
 - P. ej. los *contenidos*, los *enlaces (anchors)* o el *lenguaje*
 - Timestamp: indica la versión, siendo la más reciente la primera



¿Os recuerda a alguna aproximación vista anteriormente?

+ Middleware

Almacenamiento - Bigtable: infraestructura

- Una bigtable ocupa varios TB
 - Para su almacenamiento, se divide en *tabletas* de 100-200MB
 - Cada tableta se almacena como un conjunto de archivos en un determinado formato
- La infraestructura de Bigtable se encarga de la división/reconstrucción en tabletas
- La infraestructura de GFS se encarga del almacenamiento de las tabletas
- La infraestructura de Chubby se encarga de la monitorización de sus estados y del balanceo de carga

+ Middleware

Almacenamiento - Chubby

- Chubby es un servicio central de la infraestructura de Google para coordinación y almacenamiento, que provee:
 - **Acceso sincronizado** a actividades distribuidas
 - Bloqueos, semáforos, exclusión mutua, etc.
 - **Servicio de nombres** dentro de Google
 - **Servicio de archivos** para archivos de tamaño pequeño
 - Complementa a GFS
 - Servicio de **elección de réplicas** para lectura/escritura
- Puede parecer que no cumple el principio de simplicidad (*hacer una sola cosa bien hecha*) pero en el fondo, todas sus capacidades se derivan de un servicio central de **consenso distribuido: Paxos**

+ Middleware

Almacenamiento - Chubby

- Chubby es un servicio central de la infraestructura de Google para coordinación y almacenamiento, que provee:
 - **Acceso sincronizado** a actividades distribuidas
 - Bloqueos, semáforos, exclusión mutua, etc.
 - **Servicio de nombres** dentro de Google
 - **Servicio de archivos** para archivos de tamaño pequeño
 - Complementa a GFS
 - Servicio de **elección de réplicas** para lectura/escritura
- Puede parecer que no cumple el principio de simplicidad (*hacer una sola cosa bien hecha*) pero en el fondo, todas sus capacidades se derivan de un servicio central de **consenso distribuido: Paxos**

+ Middleware

Almacenamiento - Chubby

- Chubby comenzó como un sistema de bloqueo de archivos (conurrencia pesimista), evolucionando a una API para archivos pequeños
 - **Acceso sincronizado:** Acquire, TryAcquire, Release
 - **Servicio de archivos:**
 - Operaciones generales: Open, Close, Delete
 - Operaciones de archivo: GetStat, GetContentsAndStat, ReadDir, setContents, setACL
 - **Elección de réplicas primarias** en un sistema de réplicas primaria-respaldo.
 - La elección de la réplica primaria se hace en base a un consenso mediante el algoritmo Paxos*

* Distinto a, por ejemplo, el método del abusón o en anillo, donde la elección se hace en base a un criterio previamente definido (identificadores de proceso)

+ Middleware

Almacenamiento – Chubby/Paxos

- El algoritmo Paxos* permite consensuar operaciones en un sistema distribuido ***asíncrono***
- Requiere un coordinador
 - Se asume que el coordinador puede fallar
 - Puede haber varios coordinadores coexistiendo
 - El coordinador decidirá el valor de consenso
- Google utiliza este algoritmo para asegurar la consistencia entre réplicas de archivos
 - El valor de consenso se refiere a la actualización
 - Los procesos son réplicas que contienen el archivo

*Ver el tema de [coordinación y acuerdo](#)



Middleware

Computación - MapReduce

- La infraestructura de almacenamiento y comunicación maneja grandes conjuntos de datos distribuidos
 - Necesitamos operar con esos datos de manera distribuida
- MapReduce es un modelo de programación para aplicaciones que oculta al programador los aspectos distribuidos
 - Paralelización
 - Recuperación de fallos
 - Administración de los datos
 - Balanceo de carga



Middleware

Computación - MapReduce

- MapReduce sigue el patrón clásico en la computación en paralelo:
 - Dividir los datos (complejos) en trozos (simples)
 - Realizar el cómputo de los trozos simples (obteniendo resultados intermedios separados)
 - Unir los resultados intermedios en el resultado final
- Ejemplo: sumar todos los elementos de una matriz $10^3 \times 10^3$
 - Dividimos la matriz en 100 matrices 100×100
 - 100 procesadores calculan las sumas parciales
 - Se hace la suma total a partir de las 100 sumas parciales



Middleware

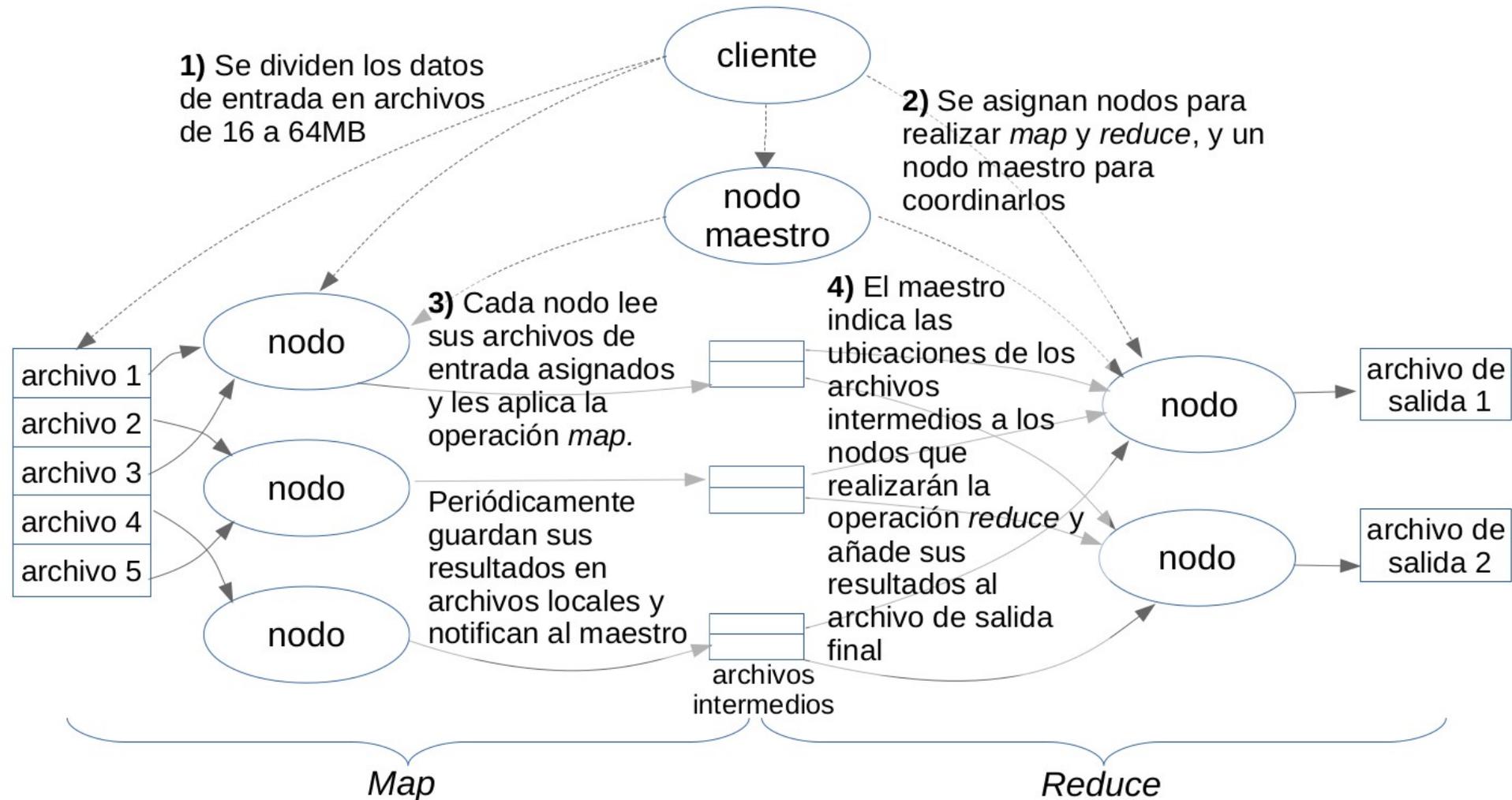
Computación - MapReduce

- MapReduce utiliza para la paralelización cuatro pasos
 - Divide los datos en trozos de igual tamaño
 - **Map**: para cada trozo realiza la operación paralelizada, produciendo uno o más pares <clave, valor>
 - Se ordenan los pares intermedios según sus claves
 - **Reduce**: se obtiene el resultado final fusionando los pares intermedios
- Ejemplo: contar el nº de veces que aparecen 1+ palabras
 - Se divide el conjunto de textos en el que buscar
 - **Map**: para cada subconjunto, cada vez que se encuentre una de las palabras buscadas, se emite un par <palabra, 1>
 - Se ordenan todos los pares según la clave (palabra)
 - **Reduce**: Se cuenta el nº de pares para cada palabra



Middleware

Computación - MapReduce



+ Middleware

Computación - MapReduce

- MapReduce se implementa en una biblioteca que oculta los detalles asociados con la paralelización y la distribución
 - Construida sobre RPC y GFS principalmente
 - Es habitual usar como entrada y salida tablas de Bigtable
 - El programador sólo tiene que centrarse en especificar las funciones *map* y *reduce*
- Este modelo de programación ofrece varias ventajas
 - Simplifica el código de los programas
 - La fase de indizado pasó de 3800 a 700 líneas de código
 - Facilita la actualización y comprensión de los algoritmos
 - Separa la lógica de la aplicación de las tareas de distribución



Middleware

Computación - Go/Sawzall

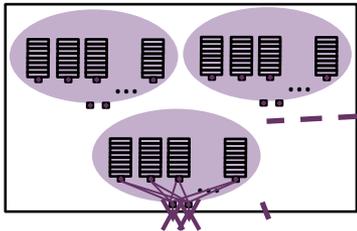
- Sawzall es un lenguaje procedural de filtros y agregadores que utiliza la filosofía de MapReduce para ofrecer operaciones más complejas (sumatorios, cuantiles, rankings, etc.)
 - Programas de 10 a 20 veces más pequeños
 - Utiliza un patrón de paralelización similar
- Sawzall muere de éxito en 2015, pues se le estaba dando uso como lenguaje de programación de propósito general.
- Se convierte en parte del sistema de logs del lenguaje de programación de propósito general Go*
- En particular, se convierte en un conjunto de bibliotecas llamado Lingo que permite seguir usando los agregadores complejos en operaciones MapReduce



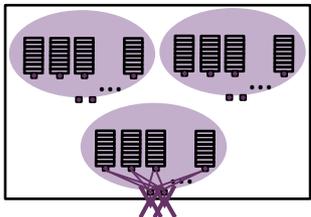
Google

Visión global

Plataforma

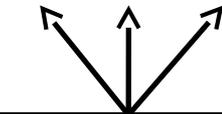


Data centers



Middleware

publish-suscribe

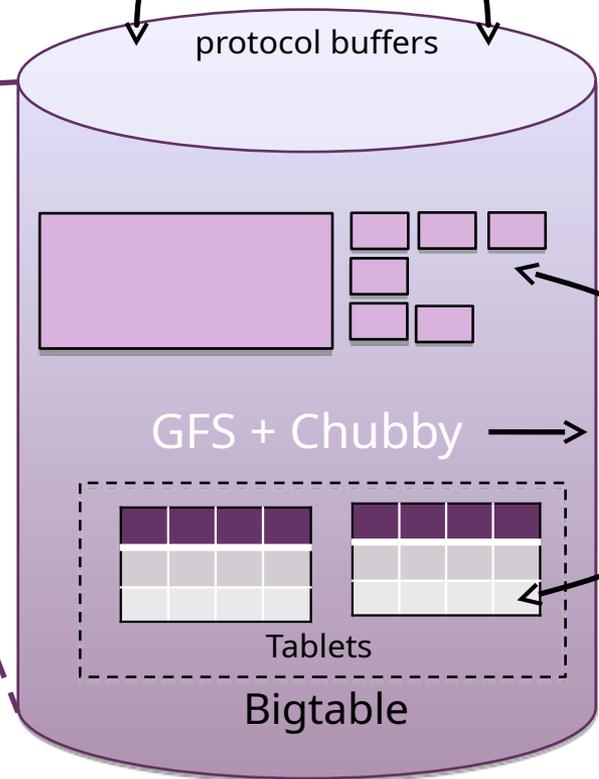
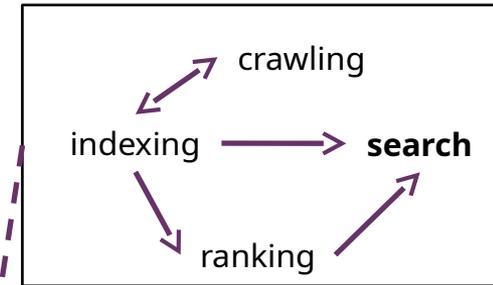


protocol buffers



Servicios

53



Google App Engine



+ Después de Google

Hadoop

- El modelo de Google, aunque no excesivamente detallado en los informes técnicos disponibles, ha dado muchas claves para el desarrollo de suites para el desarrollo de sistemas distribuidos.
- Hadoop (2006) es una de estas suites desarrollada por Apache que se inspira en Google para desarrollar sus propias soluciones:
 - Hbase como BBDD no relacional similar a Bigtable
 - HDFS como sistema de archivos similar a GFS/Colossus
 - Apache Spark como computación similar a MapReduce
- Hadoop se utiliza mucho en grandes empresas, entre ellas Yahoo, Facebook o eBay.

+ Después de Google

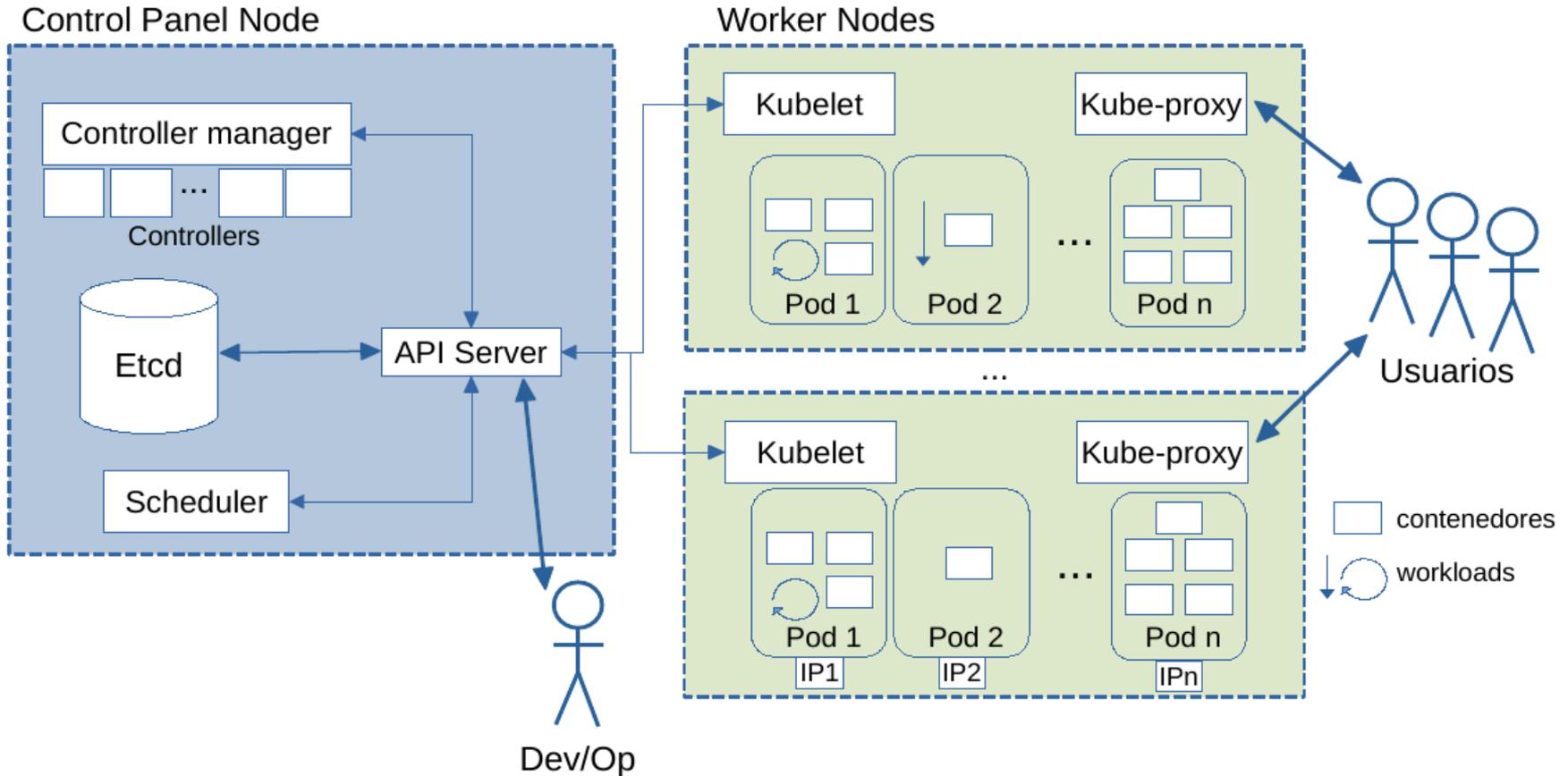
Kubernetes

- Kubernetes o K8s es un modelo más flexible de suite distribuida, liberado por Google en 2014
- Es un sistema de coordinación, en un entorno distribuido, de máquinas (físicas o virtuales) mediante contenedores
 - Un contenedor es un conjunto de programas aislado del sistema operativo, sin conocer ni depender de otros programas*.
 - Reduce los problemas de dependencias y versiones, con un sobrecoste (normalmente asumible) de memoria.
- K8s da una flexibilidad al administrador de sistemas distribuidos que le ha hecho convertirse en uno de las opciones más populares para este tipo de sistemas.

* La solución de contenedores más conocida, para el ámbito de un solo ordenador, es Docker, que permite incluir programas, o instalaciones completas de sistema operativo con sus programas, en un contenedor

+ Kubernetes

Arquitectura



* La solución de contenedores más conocida, para el ámbito de un solo ordenador, es Docker, que permite incluir programas, o instalaciones completas de sistema operativo con sus programas, en un contenedor



Kubernetes

Control Panel Node

- Un único nodo maestro toma decisiones sobre el cluster de réplicas (worker nodes).
 - Puede replicarse
- Componentes:
 - **Controladores:** procesos que gestionan distintos recursos
 - **Gestor controlador:** un proceso único que coordina a los controladores
 - **Etcad:** almacén de pares clave-valor que almacena toda la configuración del cluster y su estado.
 - **Scheduler:** proceso encargado del balanceo de carga
 - **Servidor de API:** proceso encargado de la comunicación con el desarrollador y con los otros nodos

+ Kubernetes

Worker Node

- Cada nodo en el que se ejecutan los trabajos
- Componentes:
 - **Pods:** conjunto de uno o más contenedores. Cada contenedor tiene una aplicación más sus bibliotecas y dependencias, además de una IP única.
 - **Workloads:** cada una de las tareas que se ejecutan mediante uno o más pods
 - **Kubelet:** proceso responsable de gestionar los contenedores, organizarlos en pods, y monitorizarlos.
 - **Kube-proxy:** proceso responsable del balanceo de carga y del enrutado del tráfico de los clientes.

+ Kubernetes

Workloads

- Parte del éxito de K8s son los workloads predefinidos, que permiten desplegar aplicaciones sin preocuparnos por aspectos distribuidos como la escala o la consistencia.
- Los más importantes son:
 - **Deployment:** aplicaciones sin estado. Internamente se despliegan réplicas mediante ReplicaSet, otro workload.
 - **StatefulSet:** aplicaciones con estado.
 - **DaemonSet:** aplicaciones locales de ejecución continua.
 - **Job/Cronjob:** aplicaciones locales con finalización, ejecutadas una vez/periódicamente.
- Así, un servicio estará soportado por un conjunto de tareas, que corren en pods replicados en distintos nodos de trabajo, coordinados entre sí y balanceados por el nodo de control y los kubelet/kube-proxy



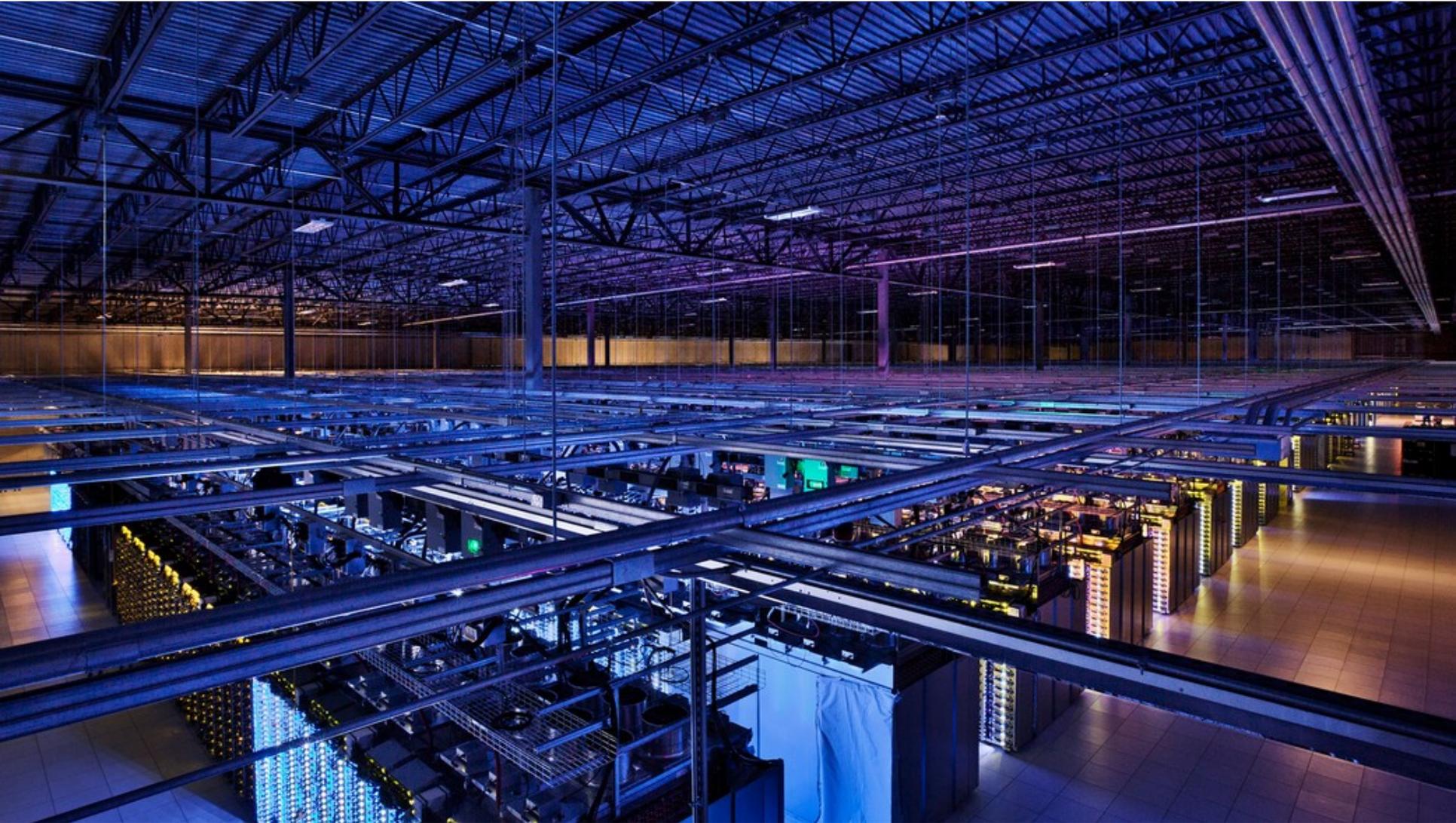
+ Google

- La **infraestructura** de **Google** es un buen **ejemplo** de casi todos los **problemas distribuidos**, y de su resolución mediante distintas aproximaciones
- Al diseñar un sistema distribuido es esencial
 - **Identificar** los **requisitos** primordiales para nuestro sistema
 - **Diseñar** una **plataforma** acorde a dichos requisitos
 - **Aprovechar** las **soluciones** distribuidas existentes o bien la **teoría** subyacente
- Google **diseña sus propias soluciones** distribuidas **basándose en soluciones conocidas**, con máxima atención en la escalabilidad y la fiabilidad.
- Para su **motor** de **búsqueda** usa tres componentes que recorren la web (**crawler**), la analizan (**indexer**) y priorizan (**ranking**)
- Para la **comunicación**, tienen una solución ad hoc basada en **RPCs** y una serialización más ligera que XML (**protocol buffers**)
- Para el **almacenamiento** distribuido, utilizan un sistema de archivos (**GFS/Colossus**) similar a NFS pero especializado en grandes archivos. La consistencia se asegura mediante **Chubby**, una versión del algoritmo Paxos de Lamport
- Para la **computación** distribuida la solución es un **map-reduce**, que se optimiza con operaciones predefinidas (**MapReduce**)
- **Hadoop** y **Kubernetes** son soluciones abiertas para construir SDD facilitando la gestión de réplicas, consistencia o balanceo de carga.



Referencias

- G. Colouris, J. Dollimore, T. Kindberg and G. Blair. *Distributed Systems: Concepts and Design (5th Ed)*. Addison-Wesley, **2011**
 - Ch. 21
- J.L. Hennessy, D.A. Patterson. *Computer Architecture, A Quantitative Approach*. 5th ed. **2012**
- Google: <http://googleblog.blogspot.com>
- S. Gherawat et al. *The Google File System*. **2003**
- Algoritmo Paxos (versión original)
 - http://en.wikipedia.org/wiki/Paxos_algorithm
- K. Marhubi. *Kubernetes from the ground up*. **2015** ([enlace](#))
- Kubernetes Documentation ([enlace](#))



Google Data Center en Council Bluffs, Iowa