

## Consistent Hashing

Supongo que deseo tener un índice distribuido sobre  $n$  máquinas/servidores  $s_1, \dots, s_n$

Possible solution: Tener una función de hash:

$$h: D \rightarrow \{1, \dots, n\}$$

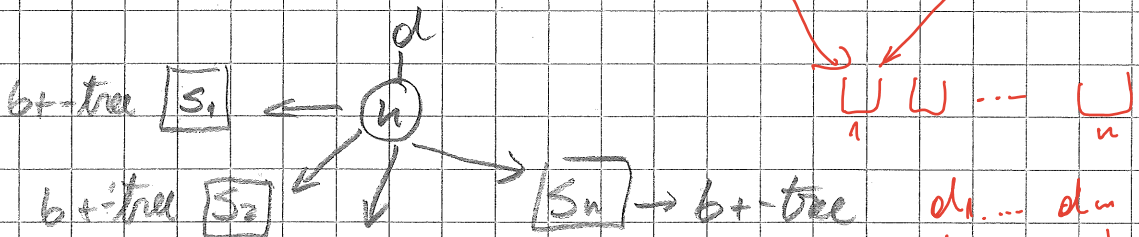
tal que cada servidor  $s_i$  recibe el conjunto

$$\{d \mid h(d) = i\}$$

Pregunta: - ¿cuál es la carga de cada máquina?

- ¿cómo buscamos/insertamos un dato?

- ¿cómo almacenamos los datos en cada servidor  $s_i$ ?



$$d_1, \dots, d_n$$

$$\frac{1}{n} \quad X_i := H \text{ de } d \text{ en } s_{i-1}$$

$$X_i = B(d, i)$$

$$\Rightarrow H(X_i) = \frac{1}{n}$$

\* ¿es esta solución más eficiente que un b+-tree en un solo servidor?

¿Cuál es el problema con esta solución?

↳ inserción de una máquina  $n+1$  requiere hacer re-hashing de todos los datos.

¿por qué quisiéramos insertar otra máquina?

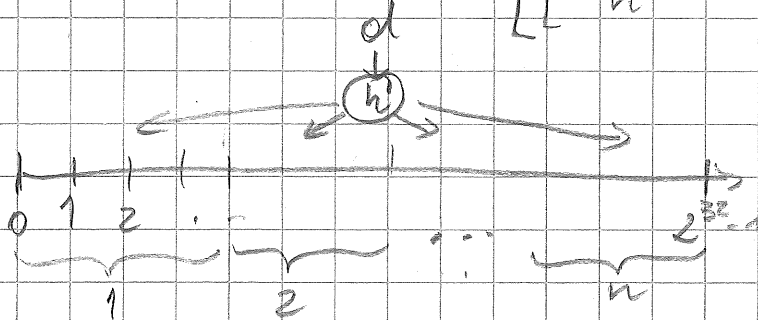
Hacer consistent hashing (otra solución)

Supongamos que en vez de tener una función de hash  $h: D \rightarrow \{1, \dots, n\}$ , tenemos

una función  $h': D \rightarrow \{0, \dots, 2^{32}-1\}$  (palabras de 32-bits)

y asignamos a cada servidor  $s_i$  el conjunto

$$\{d \mid h'(d) \in \left[ \frac{(2^{32}-1) \cdot (i-1)}{n}, \frac{(2^{32}-1) \cdot i}{n} \right) \}$$



- Preguntas:
- ¿cómo borrar/insertar un dato?
  - ¿sigue siendo "uniforme" la distribución de datos?
  - ¿es posible insertar una máquina sin hacer re-hashing de todos los datos?

### Consistent Hashing

Solución: Hacer hashing de datos y máquinas.

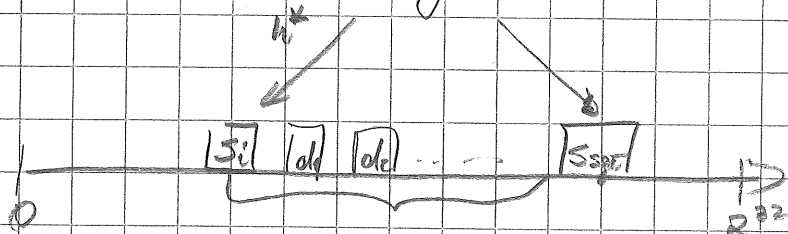
Supongamos que se tiene una función de hash:

$$h^*: D \cup N \rightarrow \{0, \dots, 2^{32}-1\}$$

y asignamos a cada servidor  $s_i$  el conjunto

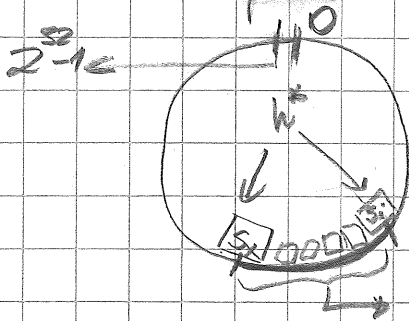
$$\{d \mid h^*(s_i) \leq h^*(d) < h^*(s_{\text{sig}}) \bmod 2^{32}\}$$

donde:  $s_{\text{sig}} = \arg \min_s \{(h^*(s) - h^*(s_i)) \bmod 2^{32}\}$



todos estos datos pertenecen a la máquina  $s_i$ .

Esto también se puede ver como un círculo



todo este cluster son asignados a  $s_i$

Preguntas: - ¿Cómo busco / inserto un dato?

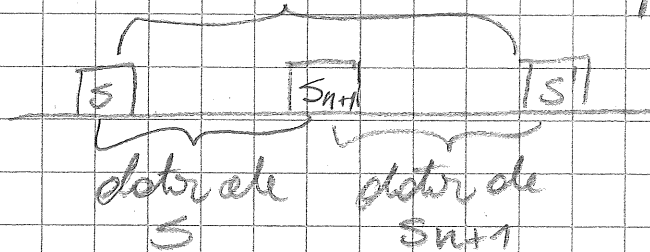
- Suposición: máquina maestra (o todos los mag.) cuentan con un b+-tree de conjuntos

$$\{ (h^*(s_i), s_i) \}$$

ordenado por  $h^*(s_i)$ .

- ¿Cómo hacemos para insertar una nueva máquina  $s_{n+1}$ ?

dato de  $S$  antes de agregar a  $s_{n+1}$



∴ - solo  $S$  debe mover sus dato de lugar

- todos los otros servidores siguen igual

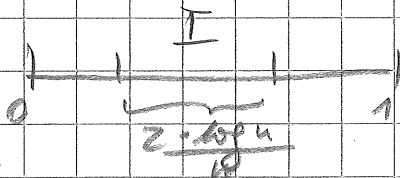
Pregunta: - ¿es uniforme la carga entre distn. max.?

Teorema: con prob.  $1 - \frac{1}{n}$ , cada máximo es responsable de a lo más  $O(\log n / n)$  por de buckets del intervalo  $[0, 2^{32} - 1]$ .

Dem.: SPDB supone que  $[0, 2^{32} - 1]$  es el intervalo  $[0, 1]$  y sea:

$$S = \{k(s_i) \mid i \in \{1, \dots, n\}\}$$

Considere el subintervalo  $I$  de tam.  $\frac{2 \cdot \log n}{n}$



$$P(I \cap S = \emptyset) = \left(1 - \frac{2 \cdot \log n}{n}\right)^n = \left(1 - \frac{2 \cdot \log n}{n}\right)^{\frac{n}{2 \cdot \log n} \cdot 2 \cdot \log n} \approx e^{-2 \cdot \log n} = \frac{1}{n^2}$$

Subdivido  $[0, 1]$  en  $I_1, \dots, I_k$  intervalos de tam.  $\frac{2 \cdot \log n}{n}$

$$P(\text{algún máx. recibe más de } \frac{4 \cdot \log n}{n} \text{ slots.}) \leq P\left(\bigcup_{i=1}^k I_i \cap S = \emptyset\right) \leq \sum_{i=1}^k P(I_i \cap S = \emptyset) = \sum_{i=1}^k \frac{1}{n^2} \leq \frac{1}{n}$$

y  $k = \left\lceil \frac{n}{2 \cdot \log n} \right\rceil$