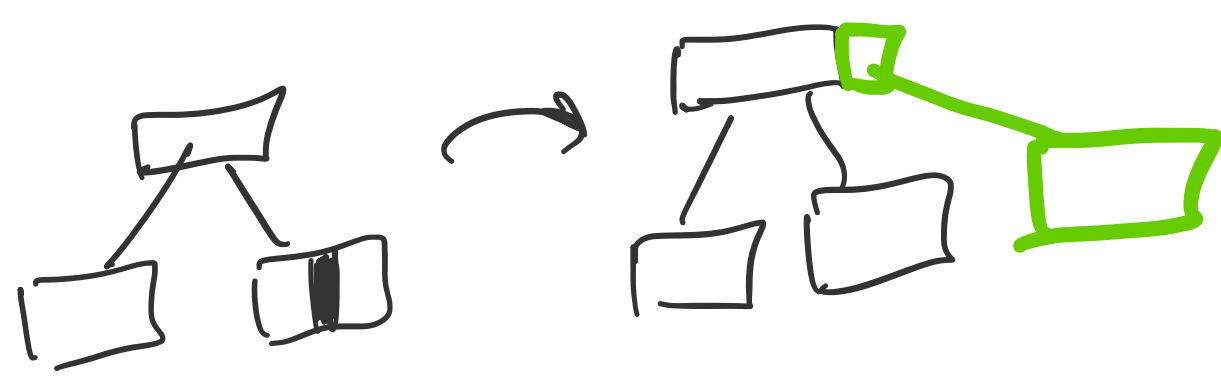
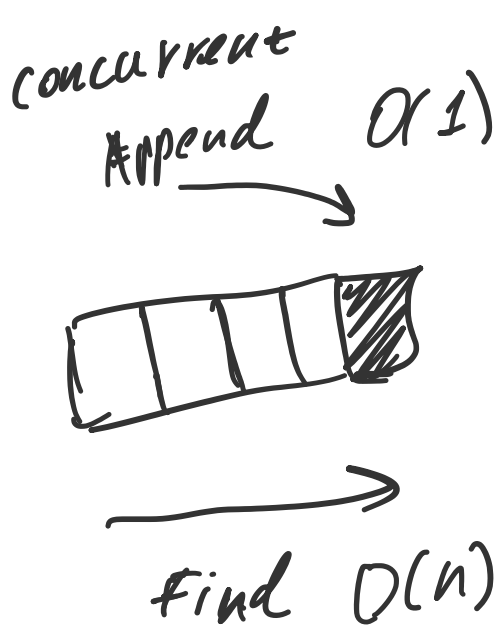


Indexing

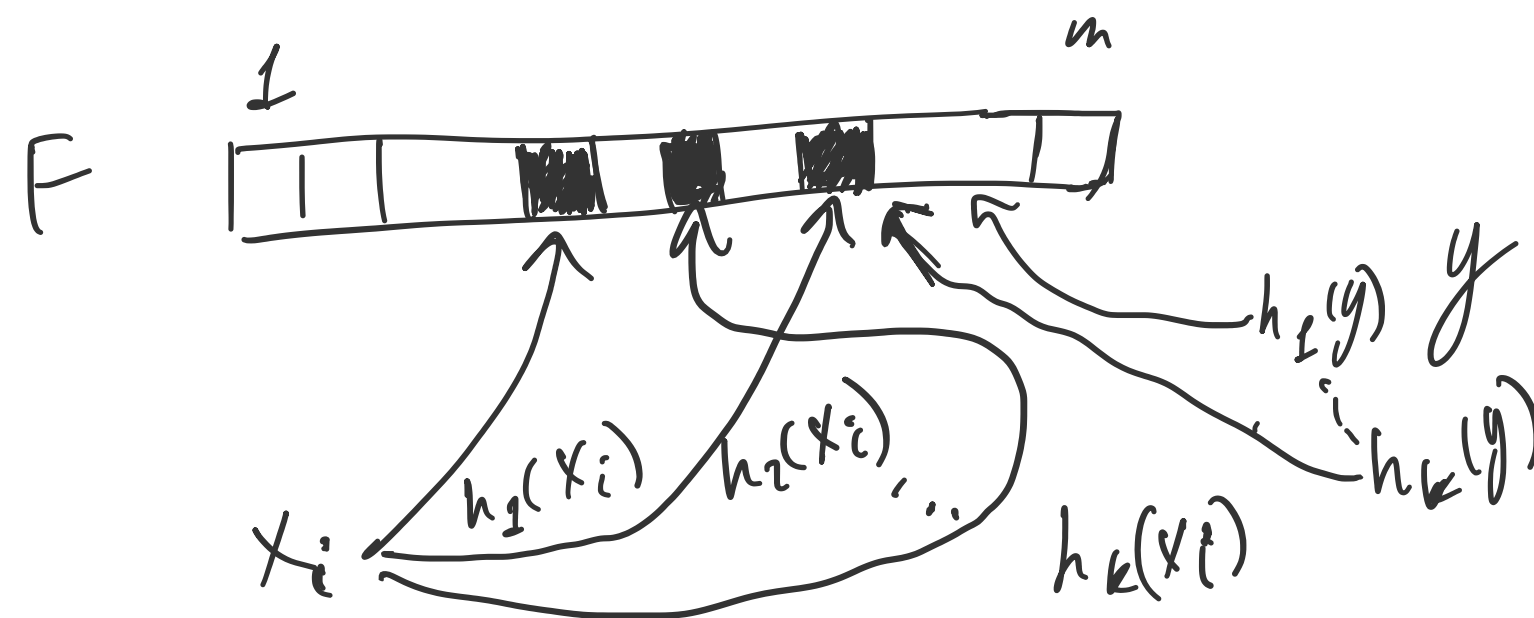
- 1) Append to file
- 2) In-memory hash index
record → offset
find O(1)
- 3) B-tree, LSM tree
B-tree fault tolerance
 - WAL
 - copy-on-write
 - concurrency latch - lock



Bloom Filter

$x_i \in K, |K|$
 $S = \{x_1, x_2, \dots, x_n\}$
 $P(x \notin S) < \epsilon$
 $\text{contain}(S, x) = \begin{cases} \text{true}, & P(x \in S) \geq 1 - \epsilon \\ \text{false}, & x \notin S \end{cases}$

S space $n \cdot \log |K|$
Filter space $\text{const} \cdot n = m$



$y \notin S$

$$P[F[h_1(y)] = 1] = 1 - P[F[h_1(y)] = 0] = 1 - \left(1 - \frac{1}{m}\right)^{nk}$$

$$P[\text{Contains}(S, y)] \approx \left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k < \epsilon$$

$(n, m) \rightarrow k \quad P \rightarrow \min$
 $(n, \epsilon) \rightarrow m$

$$\approx \left(1 - e^{-\frac{nk}{m}}\right)^k \geq \frac{1}{P} \quad f(k) = \ln(P)$$

$$f'(k) = k \cdot \frac{1}{1 - e^{-\frac{nk}{m}}} \cdot e^{-\frac{nk}{m}} \left(-\frac{n}{m}\right) + \ln(P) = 0$$

$$k \approx \lambda \cdot \frac{m}{n}$$
$$f'(\lambda) = \ln(1 - e^{-\lambda}) + \frac{e^{-\lambda}}{1 - e^{-\lambda}} \cdot \lambda = 0$$

$$\lambda = \ln 2$$
$$1/P = \left(\frac{1}{2}\right)^{\ln 2 \cdot \frac{m}{n}} = C^{\frac{m}{n}}, \quad C < 1$$

Universal hash functions
 $h \in H \quad h: K \rightarrow \{0, \dots, m-1\}$

$$\forall k, l \in K, k \neq l$$
$$\frac{\#\{h \in H \mid h(k) = h(l)\}}{|H|} \leq \frac{|H|}{m}$$
$$\frac{\#\{h \in H \mid h(k) = h(l)\}}{|H|} \leq \frac{1}{m}$$

$h \in H$

Example

p - prime

$$H_{p,m} = \{h_{a,b} \mid a \in \mathbb{Z}_p^*, b \in \mathbb{Z}_p\}$$

$1, \dots, p-1 \quad 0, \dots, p-1$

$$h_{a,b}(k) = (ak + b) \% p \times m$$

! counter
to remove
space $m \cdot \log n$

Materialized View

	key	I	II	III
time				
10	1	2	0	
16	4	3	1	
20	0	2	1	

	key	I	II	III
time				
10	1	3	3	
16	5	10	11	
20	5	12	14	

Select Count(*)

key	time	value

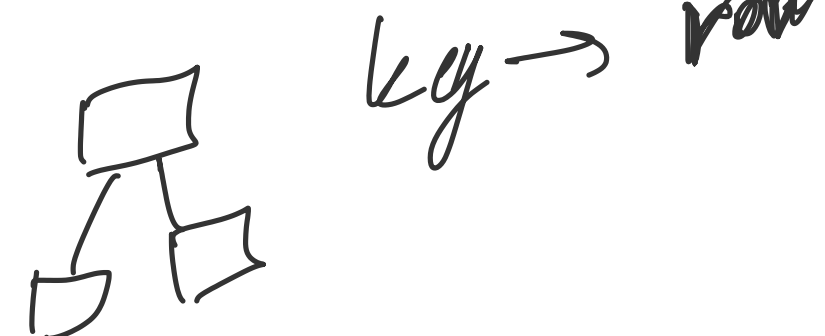
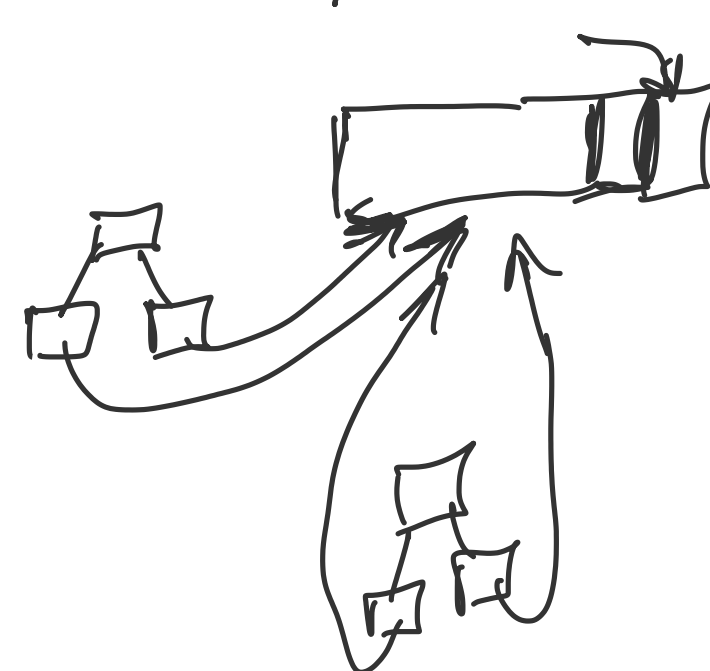
to persist data

datasync
fsync

secondary indexes

key time value

heap file vs cluster index



multicolumn index

concat VS R-tree

full-text index

query

□ □ □

document

□	□	□
□	□	□
□	□	□

Sphinx term → (doc.id, pos)

$$tf(t, d) = \frac{n_t}{\sum_k n_k}$$

$$idf(t, D) = \log \frac{|D|}{|\{d \in D \mid t \in d\}|}$$

$$tf-idf(t, d, D) = tf \cdot idf$$