

Sequential: Vector of Bits

When estimating seek costs, we need to calculate the probability distribution for the distance between two subsequent qualifying cylinders. We model the situation as a bitvector of length B with b bits set to one. Then, B corresponds to the number of cylinders and a one indicates that a cylinder qualifies.

[Later: Vector of Buckets]

Sequential: Vector of Bits (2)

Theorem

Assume a bitvector of length B . Within it b ones are uniformly distributed. The remaining $B - b$ bits are zero. Then, the probability distribution of the number j of zeros

1. between two consecutive ones,
2. before the first one, and
3. after the last one

is given by

$$\mathcal{B}_b^B(j) = \frac{\binom{B-j-1}{b-1}}{\binom{B}{b}} \quad (35)$$

Sequential: Vector of Bits (3)

Proof:

To see why the formula holds, consider the total number of bitvectors having a one in position i followed by j zeros followed by a one.

This number is

$$\binom{B-j-2}{b-2}$$

We can choose $B-j-1$ positions for i .

The total number of bitvectors is

$$\binom{B}{b}$$

and each bitvector has $b-1$ sequences of the form that a one is followed by a sequence of zeros is followed by a one.

Sequential: Vector of Bits (4)

Hence,

$$\begin{aligned} \mathcal{B}_b^B(j) &= \frac{(B-j-1)\binom{B-j-2}{b-2}}{(b-1)\binom{B}{b}} \\ &= \frac{\binom{B-j-1}{b-1}}{\binom{B}{b}} \end{aligned}$$

Part (1) follows.

To prove (2), we count the number of bitvectors that start with j zeros before the first one.

There are $B-j-1$ positions left for the remaining $b-1$ ones.

Hence, the number of these bitvectors is $\binom{B-j-1}{b-1}$ and part (2) follows.

Part (3) follows by symmetry.

Sequential: Vector of Bits (5)

We can derive a less expensive way to calculate formula for $\mathcal{B}_b^B(j)$ as follows.

For $j = 0$, we have $\mathcal{B}_b^B(0) = \frac{b}{B}$.

If $j > 0$, then

$$\begin{aligned}
 \mathcal{B}_b^B(j) &= \frac{\binom{B-j-1}{b-1}}{\binom{B}{b}} \\
 &= \frac{(B-j-1)!}{(b-1)!((B-j-1)-(b-1))!} \\
 &= \frac{B!}{b!(B-b)!} \\
 &= \frac{(B-j-1)! \ b!(B-b)!}{(b-1)!((B-j-1)-(b-1))! \ B!}
 \end{aligned}$$

Sequential: Vector of Bits (6)

$$\begin{aligned}
 \mathcal{B}_b^B(j) &= \frac{(B-j-1)! \ b!(B-b)!}{(b-1)!((B-j-1)-(b-1))! \ B!} \\
 &= b \frac{(B-j-1)! \ (B-b)!}{((B-j-1)-(b-1))! \ B!} \\
 &= b \frac{(B-j-1)! \ (B-b)!}{(B-j-b)! \ B!} \\
 &= \frac{b}{B-j} \frac{(B-j)! \ (B-b)!}{(B-b-j)! \ B!} \\
 &= \frac{b}{B-j} \prod_{i=0}^{j-1} \left(1 - \frac{b}{B-i}\right)
 \end{aligned}$$

This formula is useful when $\mathcal{B}_b^B(j)$ occurs in sums over j .

Sequential: Vector of Bits (7)

Corollary

Using the terminology of Theorem 8, the expected value for the number of zeros

1. before the first one,
2. between two successive ones, and
3. after the last one

is

$$\bar{\mathcal{B}}_b^B = \sum_{j=0}^{B-b} j \mathcal{B}_b^B(j) = \frac{B-b}{b+1} \quad (36)$$

Sequential: Vector of Bits (8)

Proof:

$$\begin{aligned}
 \sum_{j=0}^{B-b} j \binom{B-j-1}{b-1} &= \sum_{j=0}^{B-b} (B - (B-j)) \binom{B-j-1}{b-1} \\
 &= B \sum_{j=0}^{B-b} \binom{B-j-1}{b-1} - \sum_{j=0}^{B-b} (B-j) \binom{B-j-1}{b-1} \\
 &= B \sum_{j=0}^{B-b} \binom{b-1+j}{b-1} - b \sum_{j=0}^{B-b} \binom{B-j}{b} \\
 &= B \sum_{j=0}^{B-b} \binom{b-1+j}{j} - b \sum_{j=0}^{B-b} \binom{b+j}{b}
 \end{aligned}$$

Sequential: Vector of Bits (9)

$$\begin{aligned}
 \sum_{j=0}^{B-b} j \binom{B-j-1}{b-1} &= B \sum_{j=0}^{B-b} \binom{b-1+j}{j} - b \sum_{j=0}^{B-b} \binom{b+j}{b} \\
 &= B \binom{(b-1) + (B-b) + 1}{(b-1) + 1} - b \binom{b + (B-b) + 1}{b+1} \\
 &= B \binom{B}{b} - b \binom{B+1}{b+1} \\
 &= \left(B - b \frac{B+1}{b+1} \right) \binom{B}{b}
 \end{aligned}$$

With

$$\begin{aligned}
 B - b \frac{B+1}{b+1} &= \frac{B(b+1) - (Bb+b)}{b+1} \\
 &= \frac{B-b}{b+1}
 \end{aligned}$$

the claim follows

Sequential: Vector of Bits (10)

Corollary

Using the terminology of Theorem 8, the expected total number of bits from the first bit to the last one, both included, is

$$\overline{B}_{\text{tot}}(B, b) = \frac{Bb + b}{b + 1} \quad (37)$$

Sequential: Vector of Bits (11)

Proof:

We subtract from B the average expected number of zeros between the last one and the last bit:

$$\begin{aligned} B - \frac{B - b}{b + 1} &= \frac{B(b + 1)}{b + 1} - \frac{B - b}{b + 1} \\ &= \frac{Bb + B - B + b}{b + 1} \\ &= \frac{Bb + b}{b + 1} \end{aligned}$$

Sequential: Vector of Bits (12)

Corollary

Using the terminology of Theorem 8, the number of bits from the first one and the last one, both included, is

$$\overline{B}_{1\text{-span}}(B, b) = \frac{Bb - B + 2b}{b + 1} \quad (38)$$

Sequential: Vector of Bits (13)

Proof (alternative 1):

Subtract from B the number of zeros at the beginning and the end:

$$\begin{aligned}\overline{B}_{1\text{-span}}(B, b) &= B - 2\frac{B - b}{b + 1} \\ &= \frac{Bb + B - 2B + 2b}{b + 1} \\ &= \frac{Bb - B + 2b}{b + 1}\end{aligned}$$

Sequential: Vector of Bits (14)

Proof (alternative 2):

Add the number of zeros between the first and the last one and the number of ones:

$$\begin{aligned}\overline{B}_{1\text{-span}}(B, b) &= (b-1)\overline{B}_b^B + b \\ &= (b-1)\frac{B-b}{b+1} + \frac{b(b+1)}{b+1} \\ &= \frac{Bb - b^2 - B + b + b^2 + b}{b+1} \\ &= \frac{Bb - B + 2b}{b+1}\end{aligned}$$

Sequential: Applications for Bitvector Model

- If we look up one record in an array of B records and we search sequentially, how many array entries do we have to examine on average if the search is successful?
- Let a file consist of B consecutive cylinders. We search for k different keys all of which occur in the file. These k keys are distributed over b different cylinders. Of course, we can stop as soon as we have found the last key. What is the expected total distance the disk head has to travel if it is placed on the first cylinder of the file at the beginning of the search?
- Assume we have an array consisting of B different entries. We sequentially go through all entries of the array until we have found all the records for b different keys. We assume that the B entries in the array and the b keys are sorted. Further all b keys occur in the array. On the average, how many comparisons do we need to find all keys?

Sequential: Vector of Buckets

Theorem (Yao)

Consider a sequence of m buckets. For $1 \leq i \leq m$, let n_i be the number of items in a bucket i . Then there is a total of $N = \sum_{i=1}^m n_i$ items. Let $t_i = \sum_{l=0}^i n_l$ be the number of items in the first i buckets. If the buckets are searched sequentially, then the probability that j buckets that have to be examined until k distinct items have been found is

$$C_{n_i}^{N,m}(k, j) = \frac{\binom{t_j}{k} - \binom{t_{j-1}}{k}}{\binom{N}{k}} \quad (39)$$

Thus, the expected number of buckets that need to be examined in order to retrieve k distinct items is

$$\bar{C}_{n_i}^{N,m}(k) = \sum_{j=1}^m j C_{n_i}^{N,m}(k, j) = m - \frac{\sum_{j=1}^m \binom{t_{j-1}}{k}}{\binom{N}{k}} \quad (40)$$

Sequential: Vector of Buckets (2)

The following theorem is very useful for deriving estimates for average sequential accesses under different models [Especially: the above theorem follows].

Theorem (Lang/Driscoll/Jou)

Consider a sequence of N items. For a batched search of k items, the expected number of accessed items is

$$A(N, k) = N - \sum_{i=1}^{N-1} \text{Prob}[Y \leq i] \quad (41)$$

where Y is a random variable for the last item in the sequence that occurs among the k items searched.

Disk Drive Costs for N Uniform Accesses

The goal of this section is to derive estimates for the costs (time) for retrieving N cache-missed sectors of a segment S from disk.

We assume that the N sectors are read in their physical order on disk. This can be enforced by the DBMS, by the operating system's disk scheduling policy (SCAN policy), or by the disk drive controller.

Disk Drive Costs for N Uniform Accesses (2)

Remembering the description of disk drives, the total costs can be described as

$$C_{disk} = C_{cmd} + C_{seek} + C_{settle} + C_{rot} + C_{headswitch} \quad (42)$$

For brevity, we omitted the parameter N and the parameters describing the segment and the disk drive on which the segment resides.

Subsequently, we devote a (sometimes tiny) section to each summand. Before that, we have to calculate the number of qualifying cylinders, tracks, and sectors.

These numbers will be used later on.

Number of Qualifying Cylinder

- N sectors are to be retrieved.
- We want to find the number of cylinders qualifying in extent i .
- S_{sec} denotes the total number of sectors our segment contains.
- Assume: The N sectors we want to retrieve are uniformly distributed among the S_{sec} sectors of the segment.
- $S_{\text{cpe}}(i) = L_i - F_i + 1$ denotes the number of cylinders of extent i .

Disk Costs: Number of Qualifying Cylinder

The number of qualifying cylinders in extent i is:

$$S_{cpe}(i) * (1 - Prob(a\ cylinder\ does\ not\ qualify))$$

The probability that a cylinder does not qualify can be computed by deviding the total number of possibilities to chose the N sectors from sectors outside the cylinder by the total number of possibilities to chose N sectors from all S_{sec} sectors of the segment.

Hence, the number of qualifying cylinders in the considered extent is:

$$Q_C(i) = S_{cpe}(i) \mathcal{Y}_{D_{zspc}(i)}^{S_{sec}}(N) = S_{cpe}(i) \left(1 - \frac{\binom{S_{sec} - D_{zspc}(i)}{N}}{\binom{S_{sec}}{N}}\right) \quad (43)$$

Number of Qualifying Tracks

Let us also calculate the number of qualifying tracks in a partition i . It can be calculated by

$$S_{cpe}(i)D_{tpc}(1 - \text{Prob}(\text{a track does not qualify}))$$

The probability that a track does not qualify can be computed by dividing the number of ways to pick N sectors from sectors not belonging to a track divided by the number of possible ways to pick N sectors from all sectors:

$$Q_t(i) = S_{cpe}(i)D_{tpc} \mathcal{Y}_{D_{zspt}(i)}^{S_{sec}}(N) = S_{cpe}(i)D_{tpc} \left(1 - \frac{\binom{S_{sec} - D_{zspt}(i)}{N}}{\binom{S_{sec}}{N}}\right) \quad (44)$$

Number of Qualifying Tracks (2)

Just for fun, we calculate the number of qualifying sectors of an extent in zone i . It can be approximated by

$$Q_s(i) = S_{cpe}(i) D_{zspc}(i) \frac{N}{S_{sec}} \quad (45)$$

Since all $S_{cpe}(i)$ cylinders are in the same zone, they have the same number of sectors per track and we could also use Waters/Yao to approximate the number of qualifying cylinders by

$$Q_c(i) = \bar{Y}_{D_{zspc}(S_{zone}(i))}^{S_{cpe}(i) D_{zspc}(S_{zone}(i)), S_{cpe}(i)} (Q_s(i)) \quad (46)$$

If $Q_s(i)$ is not too small (e.g. > 4).

Command Costs

The command costs C_{cmd} are easy to compute. Every read of a sector requires the execution of a command. Hence

$$C_{cmd} = ND_{cmd}$$

estimates the total command costs.

Seek Costs

- often the dominant part of the costs
- we look at several alternatives from less to more precise models

Seek Costs - Upper Bound

The first cylinder we have to visit requires a random seek with cost D_{seekavg} . (Truely upper bound: $D_{\text{fseek}}(D_{\text{cyl}} - 1)$)

After that, we have to visit the remaining $Q_c(i) - 1$ qualifying cylinders.

The segment spans a total of $S_{\text{clast}}(S_{\text{ext}}) - S_{\text{cfirst}}(1) + 1$ cylinders.

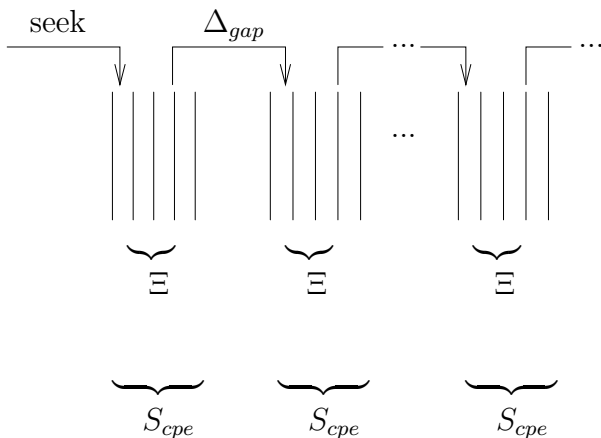
Let us assume that the first qualifying cylinder is the first cylinder and the last qualifying cylinder is the last cylinder of the segment.

Then, applying Qyang's Theorem 1 gives us the upper bound

$$C_{\text{seek}}(i) \leq (Q_c(i) - 1)D_{\text{fseek}}\left(\frac{S_{\text{clast}}(S_{\text{ext}}) - S_{\text{cfirst}}(1) + 1}{Q_c(i) - 1}\right)$$

after we have found the first qualifying cylinder.

Seek Costs - Illustration



Seek Costs - Steps

Steps:

1. Estimate for $C_{seekgap}$
2. Estimates for $C_{seekext}(i)$

Seek Costs - Interextent Costs

The average seek cost for reaching the first qualifying cylinder is D_{seekavg} . How far within the first extent are we now? We use Corollary 4 to derive that the number of non-qualifying cylinders preceding the first qualifying one in some extent i is

$$\bar{B}_{Q_c(i)}^{S_{\text{cpe}}(i)} = \frac{S_{\text{cpe}}(i) - Q_c(i)}{Q_c(i) + 1}.$$

The same is found for the number of non-qualifying cylinders following the last qualifying cylinder. Hence, for every gap between the last and the first qualifying cylinder of two extents i and $i + 1$, the disk arm has to travel the distance

$$\Delta_{\text{gap}}(i) := \bar{B}_{Q_c(i)}^{S_{\text{cpe}}(i)} + S_{\text{cfirst}}(i + 1) - S_{\text{clast}}(i) - 1 + \bar{B}_{Q_c(i+1)}^{S_{\text{cpe}}(i+1)}$$

Using this, we get

$$C_{\text{seekgap}} = D_{\text{seekavg}} + \sum_{i=1}^{S_{\text{ext}}-1} D_{\text{fseek}}(\Delta_{\text{gap}}(i))$$

Seek Costs - Intraextent Costs (2)

Let us turn to $C_{seekext}(i)$. We first need the number of cylinders between the first and the last qualifying cylinder, both included, in extent i . It can be calculated using Corollary 6:

$$\Xi_{ext}(i) = \overline{B}_{1-span}(S_{cpe}(i), Q_c(i))$$

Hence, $\Xi(i)$ is the minimal span of an extent that contains all qualifying cylinders.

Seek Costs - Intraextent Costs

Using $\Xi(i)$ and Qyang's Theorem 1, we can derive an upper bound for $C_{seekext}(i)$:

$$C_{seekext}(i) \leq (Q_c(i) - 1)D_{fseek} \left(\frac{\Xi(i)}{Q_c(i) - 1} \right) \quad (47)$$

Alternatively, we could formulate this as

$$C_{seekext}(i) \leq (Q_c(i) - 1)D_{fseek} (\bar{B}_{Q_c(i)}^{S_{cpe}(i)}) \quad (48)$$

by applying Corollary 4.

Seek Costs - Intraextent Costs (2)

A seemingly more precise estimate for the expected seek cost within the qualifying cylinders of an extent is derived by using Theorem 8:

$$C_{seekext}(i) = Q_c(i) \sum_{j=0}^{S_{cpe}(i) - Q_c(i)} D_{fseek}(j+1) \mathcal{B}_{Q_c(i)}^{S_{cpe}(i)}(j) \quad (49)$$

Settle Costs

The average settle cost is easy to calculate. For every qualifying cylinder, one head settlement takes place:

$$C_{settle}(i) = Q_c(i)D_{rdsettle} \quad (50)$$

Rotational Delay

Let us turn to the rotational delay.

For some given track in zone i ,

we want to read the $Q_t(i)$ qualifying sectors contained in it.

On average, we would expect that the read head is ready to start reading in the middle of some sector of a track.

If so, we have to wait for $\frac{1}{2}D_{zscan}(S_{zone}(i))$ before the first whole sector occurs under the read head.

However, due to track and cylinder skew, this event does not occur after a head switch or a cylinder switch.

Instead of being overly precise here, we ignore this half sector pass by time and assume we are always at the beginning of a sector.

This is also justified by the fact that we model the head switch time explicitly.

Rotational Delay (2)

Assume that the head is ready to read at the beginning of some sector of some track.

Then, in front of us is a — cyclic, which does not matter — bitvector of qualifying and non-qualifying sectors.

We can use Corollary 5 to estimate the total number of qualifying and non-qualifying sectors that have to pass under the head until all qualifying sectors have been seen.

The total rotational delay for the tracks of zone i is

$$C_{rot}(i) = Q_t(i) D_{zscan}(S_{zone}(i)) \bar{B}_{tot}(D_{zspt}(S_{zone}(i)), Q_{spt}(i))$$

where $Q_{spt}(i) = \overline{\mathcal{W}}_1^{S_{sec}, D_{zspt}(S_{zone}(i))}(N) = D_{zspt}(S_{zone}(i)) \frac{N}{S_{sec}}$ is the expected number of qualifying sectors per track in extent i . In case $Q_{spt}(i) < 1$, we set $Q_{spt}(i) := 1$.

Rotational Delay (3)

A more precise model is derived as follows.

We sum up for all j the product of (1) the probability that j sectors in a track qualify and (2) the average number of sectors that have to be read if j sectors qualify.

This gives us the number of sectors that have to pass the head in order to read all qualifying sectors.

We only need to multiply this number by the time to scan a single sector and the number of qualifying tracks.

We can estimate (1) using Theorem 7. For (2) we again use Corollary 5.

$$C_{rot}(i) = S_{cpe}(i) D_{tpc} D_{zscan}(S_{zone}(i)) \\ * \sum_{j=1}^{\min(N, D_{zspt}(S_{zone}(i)))} \chi_{D_{zspt}(S_{zone}(i))}^{S_{sec}}(N, j) \bar{B}_{tot}(D_{zspt}(S_{zone}(i)), j)$$

Rotational Delay (4)

Yet another approach:

Split the total rotational delay into two components:

1. $C_{rotpass}(i)$ measures the time needed to skip unqualifying sectors
2. $C_{rotread}(i)$ that for scanning the qualifying sectors

Then

$$C_{rot} = \sum_{i=1}^{S_{ext}} C_{rotpass}(i) + C_{rotread}(i)$$

where the total transfer cost of the qualifying sectors can be estimated as

$$C_{rotread}(i) = Q_s(i) D_{zscan}(S_{zone}(i))$$

Rotational Delay (5)

Let us treat the first component ($C_{rotpass}(i)$).

Assume that j sectors of a track in extent i qualify.

The expected position on a track where the head is ready to read is the middle between two qualifying sectors.

Since the expected number of sectors between two qualifying sectors is $D_{zspt}(S_{zone}(i))/j$, the expected number of sectors scanned before the first qualifying sector comes under the head is

$$\frac{D_{zspt}(S_{zone}(i))}{2j}$$

Rotational Delay (6)

The expected positions of j qualifying sectors on the same track is such that the number non-qualifying sectors between two successively qualifying sectors is the same.

Hence, after having read a qualifying sector $\frac{D_{zspt}(S_{zone}(i))}{j}$ unqualifying sectors must be passed until the next qualifying sector shows up.

The total number of unqualifying sectors to be passed if j sectors qualify in a track of zone i is

$$N_s(j, i) = \frac{D_{zspt}(S_{zone}(i))}{2j} + (j - 1) \frac{D_{zspt}(S_{zone}(i)) - j}{j}$$

Rotational Delay (7)

Using again Theorem 7, the expected rotational delay for the unqualifying sectors then is

$$C_{rotpass}(i) = S_{cpe}(i) D_{tpc} D_{zscan}(S_{zone}(i)) \\ \min(N, D_{zspt}(S_{zone}(i))) \\ * \sum_{j=1}^{D_{zspt}(S_{zone}(i))} \mathcal{X}_{D_{zspt}(S_{zone}(i))}^{S_{sec}}(N, j) N_s(j, i)$$

Head Switch Costs

The average head switch cost is equal to the average number of head switches that occur times the average head switch cost.

The average number of head switch is equal to the number of tracks that qualify minus the number of cylinders that qualify since a head switch does not occur for the first track of each cylinder.

Summarizing

$$C_{headswitch} = \sum_{i=1}^{S_{ext}} (Q_t(i) - Q_c(i)) D_{hdswitch} \quad (51)$$

where Q_t is the average number of tracks qualifying in an extent.

Discussion

We neglected many problems in our disk access model:

- partially filled cylinders,
- pages larger than a block,
- disk drive's cache,
- remapping of bad blocks,
- non-uniformly distributed accesses,
- clusteredness,
- and so on.

Whereas the first two items are easy to fix, the rest is not so easy.