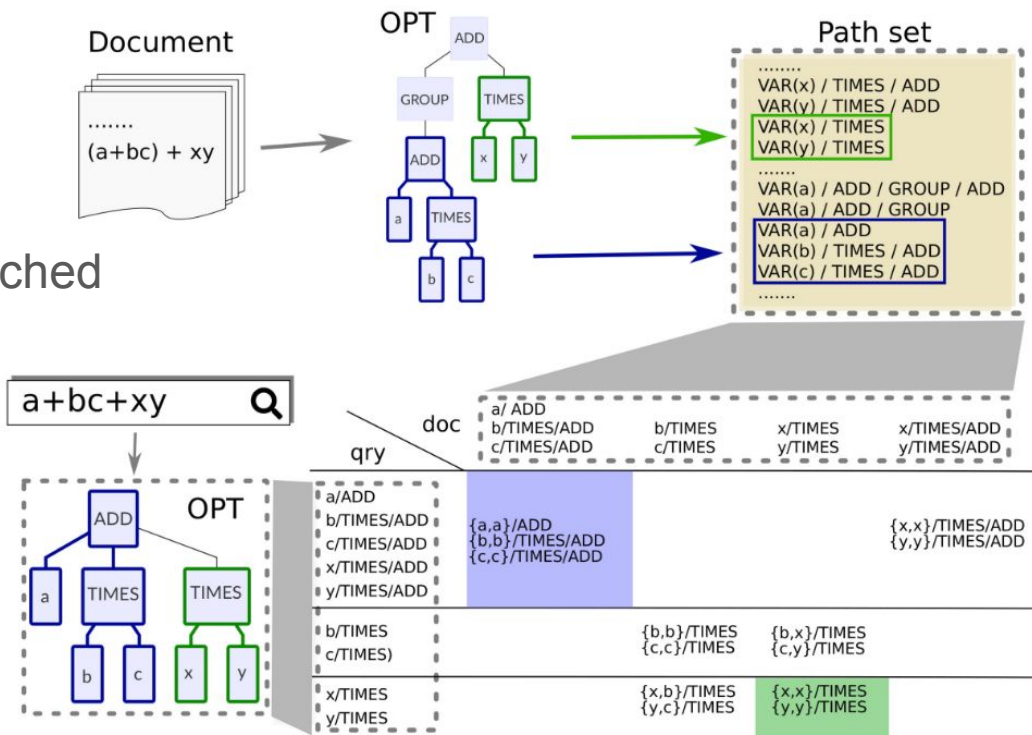


Section tree pruning for finding best matched subexpression

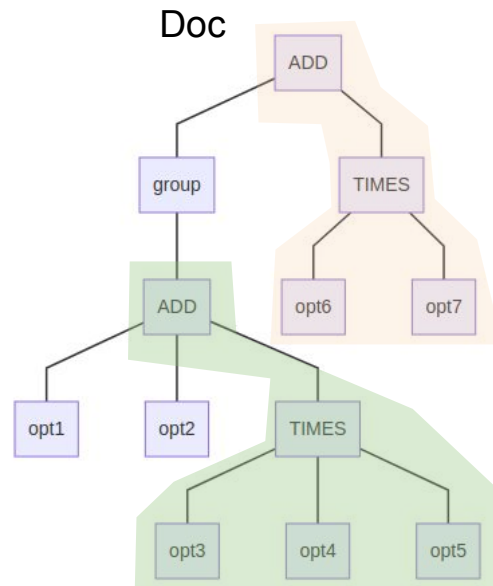
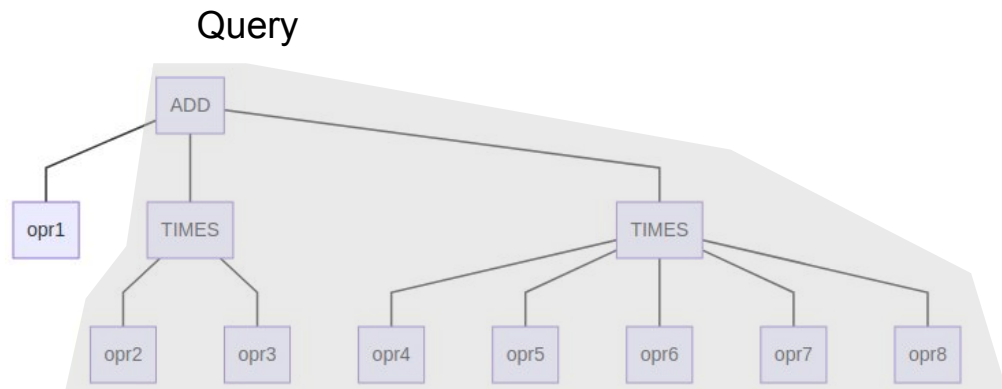
Wei Zhong

Recap

- We want to find the best-matched subexpressions.
- The experiment shows the model is computation-bounded.
- In the alignment algo, we need to examine all rows of table and all K best-matched subexpressions.
- Pruning idea: Each tree rooted at some node has upper-bound of matched nodes and we can skip rows with upperbound \leq a threshold.



Sector tree



Sector tree: Subtrees with same leaf-root paths.

The maximum matching between two sector trees $T1$ and $T2$ is just $\min(\text{leaves}(T1), \text{leaves}(T2))$. Old algorithm tries to match tree paths one by one, if we index sector tree (no leaf IDs under it), it will speed up the matching algorithm.

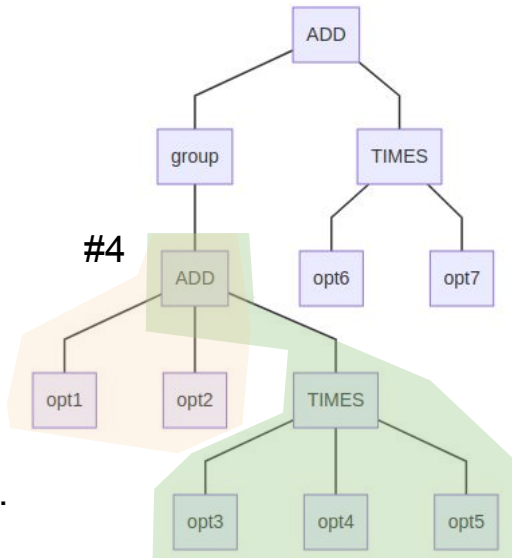
Sector tree

One rooted subtree can be decomposed into several sector trees.

Each sector tree is stored in only one posting list, and it is impossible for one posting list sector trees to match any paths of sector trees from other posting lists.

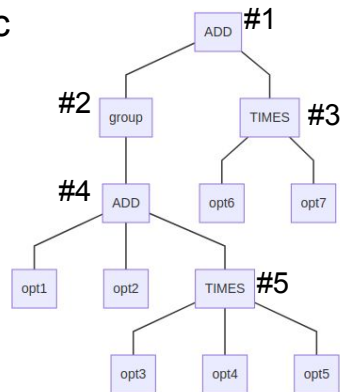
For matching, a pair of value is enough to store a sector subtree. One is the node ID of the subtree root-end node, the other is the “width” of the subtree. Use the notation **(rootID, width)** for sector tree.

For example, the subtree rooted at node#12 has two sector trees: **(4, 2)** and **(4, 3)**.

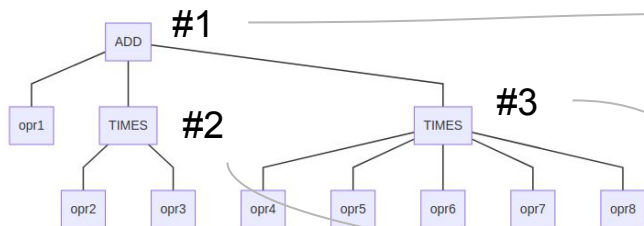


Sector tree posting lists

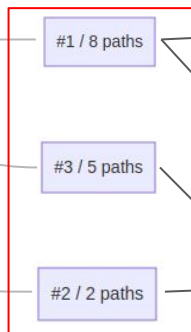
Doc



Query OPT

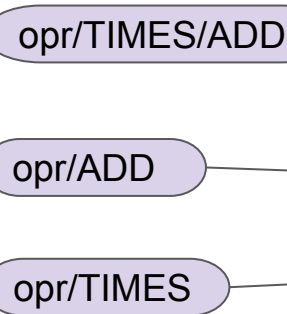


Query nodes

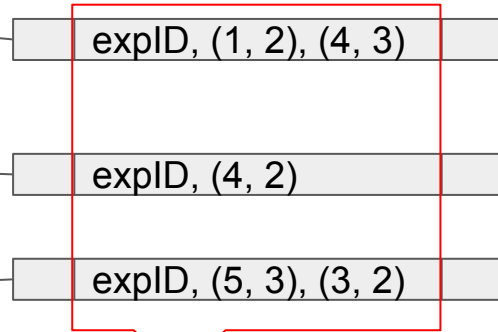


Ordered by upper-bound value

Query paths



Doc posting lists

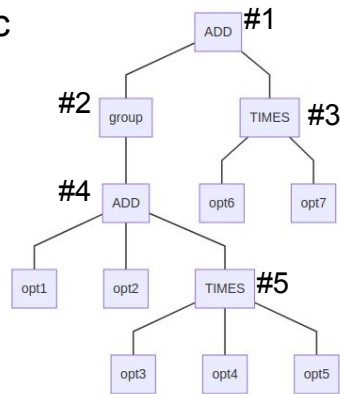


A single doc expression at a time.

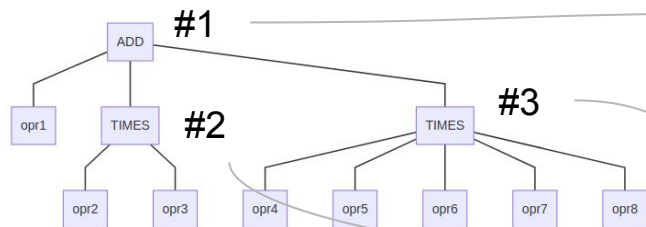
Construct sector posting lists, and utilize query node score upper-bounds.

Sector-tree width values

Doc



Query OPT



Query nodes

#1 / 8 paths

#3 / 5 paths

#2 / 2 paths

Query paths

opr/TIMES/ADD

opr/ADD

opr/TIMES

Doc posting lists

expID, (1, 2), (4, 3)

[2, 0, 0, 3, 0]

expID, (4, 2)

[0, 0, 0, 2, 0]

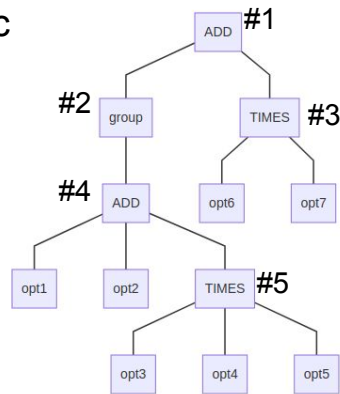
expID, (5, 3), (3, 2)

[0, 0, 2, 0, 3]

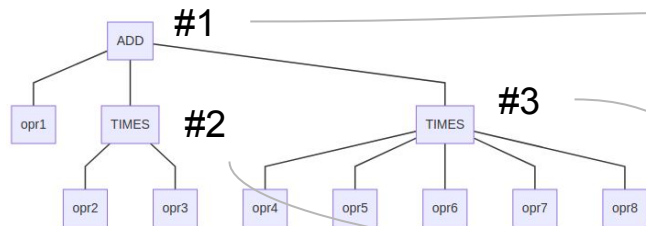
Assign *sector-tree width* values to query node -> paths relation, and use vector representation in doc space.

Sector-tree matching iterations

Doc



Query OPT



Query nodes

#1 / 8 paths

#3 / 5 paths

#2 / 2 paths

Query paths

opr/TIMES/ADD

opr/ADD

opr/TIMES

Doc posting lists

explD, (1, 2), (4, 3)

[2, 0, 0, 3, 0]

explD, (4, 2)

[0, 0, 0, 2, 0]

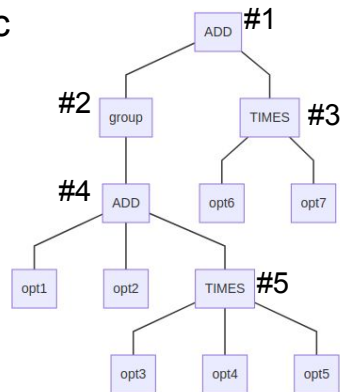
explD, (5, 3), (3, 2)

[0, 0, 2, 0, 3]

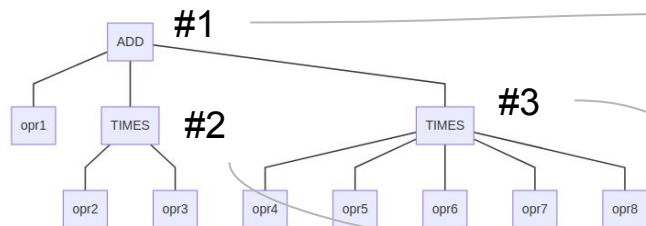
Iteration (1): $\min(7, [2, 0, 0, 3, 0])$, unmatched query paths: $7 - \max(2, 3) = 4$ (dynamic upperbound $\rightarrow 8 - 4 = 4$)
 $+ \min(1, [0, 0, 0, 2, 0])$, unmatched query paths: $1 - \max(1) = 0$
 $= [2, 0, 0, 4, 0]$, where the widest = 4.

Sector-tree matching iterations

Doc



Query OPT



Query nodes

#1 / 8 paths

#3 / 5 paths

#2 / 2 paths

Query paths

opr/TIMES/ADD

opr/ADD

opr/TIMES

Doc posting lists

expID, (1, 2), (4, 3)
[2, 0, 0, 3, 0]

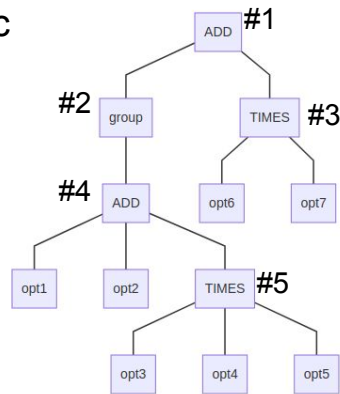
expID, (4, 2)
[0, 0, 0, 2, 0]

expID, (5, 3), (3, 2)
[0, 0, 2, 0, 3]

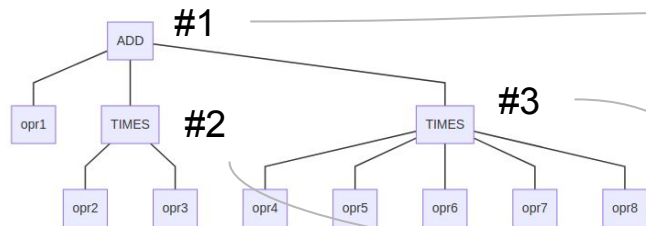
Iteration (2): $\min(5, [0, 0, 2, 0, 3])$, unmatched query paths: $5 - \max(2, 3) = 2$
 = $[0, 0, 2, 0, 3]$, update the current widest to $\max(\text{previous widest}, 3) = 4$ (not changed)

Sector-tree pruning (1)

Doc



Query OPT



Query nodes

#1 / 8 paths

#3 / 5 paths

#2 / 2 paths

Query paths

opr/TIMES/ADD

opr/ADD

opr/TIMES

Doc posting lists

expID, (1, 2), (4, 3)
[2, 0, 0, 3, 0]

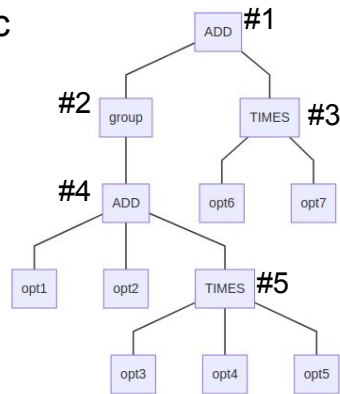
expID, (4, 2)
[0, 0, 0, 2, 0]

expID, (5, 3), (3, 2)
[0, 0, 2, 0, 3]

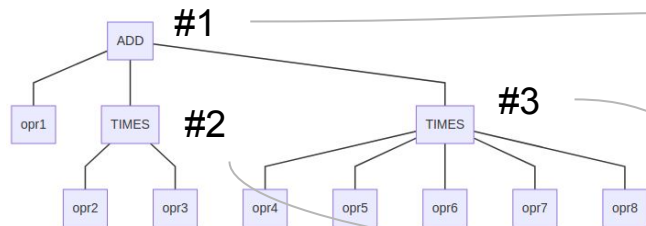
Iteration (3): Since this node has upper-bound value of $2 \leq$ current widest value 4, we can skip this iteration.

Sector-tree pruning (2)

Doc



Query OPT



Query nodes

#1 / 8 paths

#3 / 5 paths

#2 / 2 paths

Query paths

opr/TIMES/ADD

opr/ADD

opr/TIMES

Doc posting lists

expID, (1, 2), (4, 3)
[2, 0, 0, 3, 0]

expID, (4, 2)
[0, 0, 0, 2, 0]

expID, (5, 3), (3, 2)
[0, 0, 2, 0, 3]

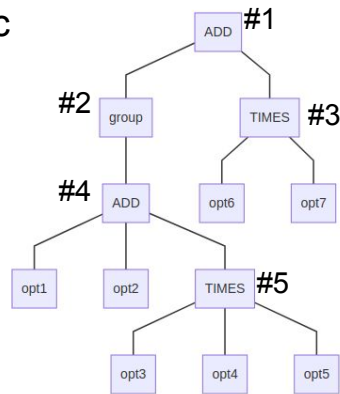
Go back to iteration (1):

From min(7, [2, 0, 0, 3, 0]), we know unmatched query paths: $7 - \max(2, 3) = 4$ (dynamic upper bound $\rightarrow 8 - 4 = 4$)

If the dynamic upper bound value \leq current widest value, we can also skip that iteration as early as possible.

Sector-tree pruning (3)

Doc

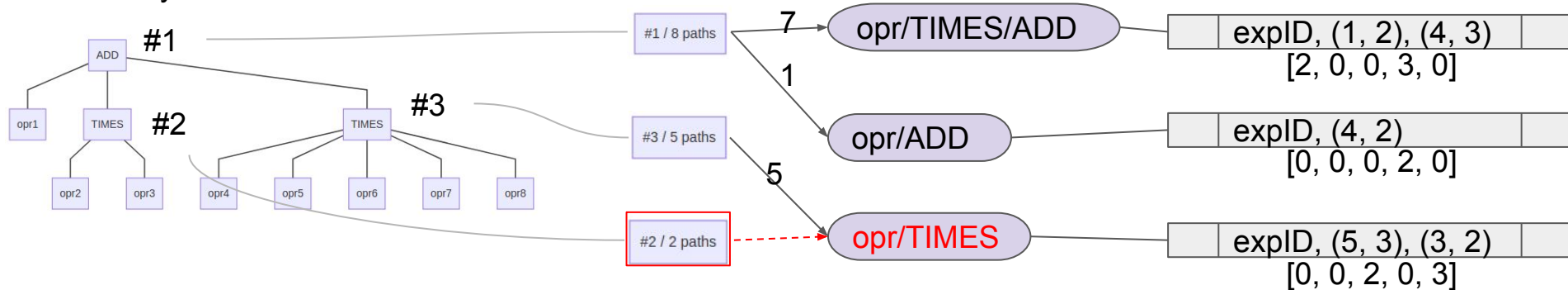


Query OPT

Query nodes

Query paths

Doc posting lists



Also, assume we have a min-heap that stores currently top-K best scored hits, at the top of the min-heap is the lowest scored top-K hit, denote the matched leaves of it as threshold value " Θ ".

Then if *upper bound* value $\leq \Theta$, we can decrease the reference counter of associated posting list, and if this posting list reference counter goes ≤ 0 , that posting list can be dropped completely for later merging.

Pseudo code

```
1  struct sector_tr {
2      int rnode;
3      int width;
4  };
5
6  struct posting {
7      struct sector_tr *(*get_sector_trees());
8      int (*get_min());
9      void (*advance());
10     int expID;
11     int refcnt;
12 };
13
14 struct pruner_node {
15     int width; // each node is ordered by this key
16     struct sector_tr secttr[64]; // each sect rnode is the same here
17     struct posting *posting[64];
18 };
19
20 int main() // simulate one merge stage iteration
21 {
22     struct posting all_post[64];
23     struct pruner_node q_nodes[64];
24     int theta = 123; /* best matched number of leaves */
25     int q_widest_match = 0;
26     for (int i = 0; i < 64 /* |T_q| */; i++) {
27         struct pruner_node *q_node = q_nodes + i;
28         int q_upperbound = q_node->width;
29         int q_node_match = 0;
30         int sumvec[128] = {0};
31
32         if (q_upperbound <= theta) {
33             for (int j = 0; j < 64 /* |l(T_q)| */; j++) {
34                 struct posting *post = q_node->posting[j]
35                 post->refcnt --;
36             }
37             /* delete this pruner_node */;
38             continue;
39         }
40
41         for (int j = 0; j < 64 /* |l(T_q)| */; j++) {
42             struct sector_tr *q_secttr = q_node->secttr + j;
43             struct posting *post = q_node->posting[j];
44             struct sector_tr *d_secttr = post->get_sector_trees();
45             int qw = q_secttr->width;
46             int max = 0;
47
48             /* skip non-hit postings */
49             if (post->get_min() != post->expID)
50                 continue;
51
52             /* vector addition */
53             for (int k = 0; k < 64 /* |l(T_d)| */; k++) {
54                 int dw = d_secttr[k].width;
55                 int rn = d_secttr[k].rnode;
56                 if (qw > dw) {
57                     sumvec[rn] += dw;
58                     if (dw > max) max = dw;
59                 } else {
60                     sumvec[rn] += qw;
61                     if (qw > max) max = qw;
62                 }
63             }
64
65             if (sumvec[rn] > q_node_match) {
66                 q_node_match = sumvec[rn];
67             }
68
69             /* update upperbound for each vector addition */
70             if (max < qw) {
71                 q_upperbound -= (qw - max);
72                 if (q_upperbound <= theta)
73                     break;
74             }
75         }
76
77         if (q_node_match > q_widest_match)
78             q_widest_match = q_node_match;
79     }
80 }
```