

Assignment Project Exam Help

Add WeChat edu\_assist\_pro

# Advanced Safe DAAT

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Wei Wang

UNSW

# Assignment Project Exam Help

## DAAT Add WeChat edu\_assist\_pro

- Idea

- Access and score each document before moving to the next, based on the inverted index

- Invariant

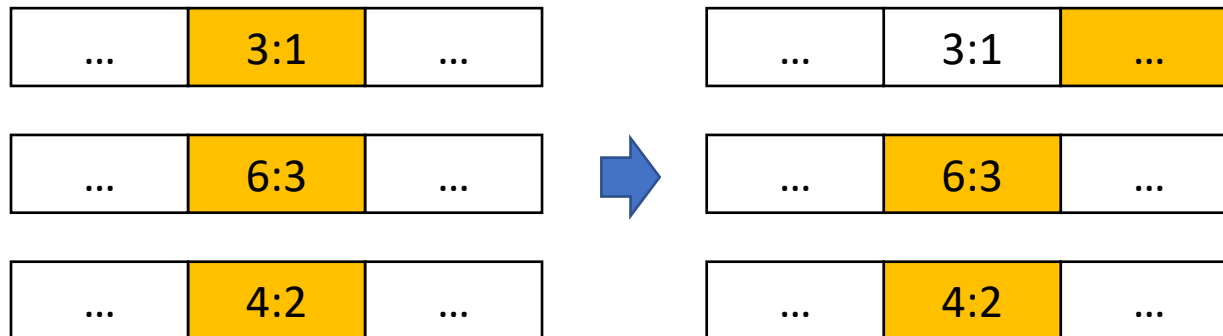
- All the documents with  $curDoc$  has been processed

<https://eduassistpro.github.io/>  
#doc to lists of the query

Add WeChat edu\_assist\_pro

Always move the smaller

... : curDoc



# Assignment Project Exam Help

DAAT Add WeChat edu\_assist\_pro

- Top-k optimization
  - Current top-k-th document's score = threshold
- Optimization
  - No need to access/score documents if  $e \leq \text{threshold}$ 
    - ➔ Skipping docIDs, but how?
- Preliminaries:
  - $UB(w)$ : upper bound of the score contribution of any document in  $w$ 's postings list

$$UB(w) = \text{idf}(w) \cdot \max_{e \in L} e.\text{tf}$$



← Assume no normalization on the raw tf,  $UB(w) = \text{idf}(w) * 4$

# Assignment Project Exam Help

## Idea 1 Add WeChat edu\_assist\_pro

- Consider  $Q = A \ B \ C \ D$

- Assume  $|A| \leq |B| \leq |C| \leq |D|$

- Can we split the query into, e.g.,  $Q1 = A \ OR \ B$ , and  $Q2 = C \ D$

- Answering  $Q1$  is more for the final scoring

fully only accessing lists in  $Q2$

<https://eduassistpro.github.io/>

- Working out a sufficient condition for the above

Add WeChat edu\_assist\_pro

Generating the candidates

Scoring the candidates

A	<table><tr><td>...</td><td>*.*</td><td>...</td></tr></table>	...	*.*	...	UB = 5
...	*.*	...			
B	<table><tr><td>...</td><td>*.*</td><td>...</td></tr></table>	...	*.*	...	UB = 3
...	*.*	...			
C	<table><tr><td>...</td><td>3:*</td><td>...</td></tr></table>	...	3:*	...	UB = 2
...	3:*	...			

MaxScore(d3) = ?, if d3 does not exist in A or B's list

➔ \_\_\_\_

If threshold  $\geq 2$ , what can you infer?

## Assignment Project Exam Help

# Determine the actual Terms

- Only need to focus on documents that occurred **ONLY** in the lists of optional terms → Estimate their upper bounding score
- Algorithm:
  - sort terms in decreasing order of their cumulative UB values
  - find the largest suffix of the terms such that the cumulative UB values is larger than the threshold (current top- $n$ 's score)

# Assignment Project Exam Help

## Simplified MaxS

- Assume all idf = 1, threshold = 2, and  $Q1 = A \ B$  (required terms)

- Step 1: generate candi

- $T = \text{Result}(Q1)$

- Step 2: score candidates and get to

- Foreach d in T: optional terms
    - $\text{Score}(d) \text{ /* using lists in } Q2 = C \text{ */}$
  - Keep top-k documents as the answer

Any problem with this alg?

- Large candidate size:
  - $[|A|, |A|+|B|]$
- The same threshold is used throughout



Killing two birds  
with one stone!

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	...	1:3	...	UB = 5	$T = \{1, 23\}$	C	1:1	33:1	55:1
B	...	23:2	...	UB = 3					

$\text{score}(d1) = 4$   
 $\text{score}(d23) = 2$

# Assignment Project Exam Help

## MaxScore

- [Step 0] Obtain the initial threshold  $\alpha$

- Update the required terms

- Repeat until the *stop*

[Step 1] • Perform one DAAT step  $\text{score}(d, s1)$ , where  $s1$  is the partial score of  $d$  on required terms

[Step 2] • Score  $d$  using the optional terms via  $s_k$

[maintenance] • If  $d$ 's final score is larger than  $\alpha$

- Update the top-k results
  - Update  $\alpha$
  - Update the required terms and optional terms

- Fixed order of terms (decreasing UB values)
    - From the rear of the list, find the maximum number of terms such that their UB sum  $\leq \alpha$

## Example (k=2)

- Step 0:
  - $d1 = 11, d2 = 7$
  - $\alpha =$
  - $Q1 =$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>



# Assignment Project Exam Help

Fixed order = C B A

## Example (k=2)

- $\alpha = 7$
- Iter 1:
  - curDoc = d5
  - partial score = 1
  - “probe” A to get full score = 1
  - Nothing to update

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

# Assignment Project Exam Help

Fixed order = C B A

## Example (k=2)

- $\alpha = 7$
- Iter 2:
  - curDoc = d7
  - partial score = 9
  - “probe” A to get full score = 10
  - Update
    - $\alpha = 10$
    - Q1 = C

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

# Assignment Project Exam Help

Fixed order = C B A

## Example (k=2)

- $\alpha = 10$
- Iter 3:
  - curDoc = d10
  - partial score = 1
  - “probe” A and B to get full score = 1
  - Nothing to update

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

# Assignment Project Exam Help

Fixed order = C B A

## Example (k=2)

- $\alpha = 10$
- Iter 3:
  - curDoc = d11
  - partial score = 8
  - “probe” A and B to get full score = 13
  - Update
    - $\alpha = 11$
    - Q1 = C
- End

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

A	B	C
UB <sub>A</sub> = 4	UB <sub>B</sub> = 5	UB <sub>C</sub> = 8
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

- This is the typical stopping condition; There is another possible stopping condition though. Can you figure it out? (not in this example)
- We can further optimize the algorithm to remove unnecessary “probes” on Q2 list(s). Can you find it out?

# Assignment Project Exam Help

## Idea 2 Add WeChat edu\_assist\_pro

Sorted Term	A	C	B
Doc	*	*	*
Cumulative Upper Bound	5	7	10

A	...	2:*	...	UB = 5
B	...	7:*	...	UB = 3
C	...	4:*	...	UB = 2

if $\alpha = 9$ , the first document that can score above $\alpha$ is from __B__
Pivot document is the current d <a href="https://eduassistpro.github.io/">https://eduassistpro.github.io/</a>

Pivot term = B

Pivot doc = d7

Is it possible for any doc < PivotDoc to enter int

Add WeChat edu\_assist\_pro

Case I: smallest DID  $\neq$  PivotDoc

Sorted Term	A	C	B
Doc	2	*	7
Cumulative Upper Bound	5	7	10

score(d2)  $\leq$  8

→ align preceding lists to PivotDoc:

A.skipTo(d7), C.skipTo(d7)

→ Check again

# Assignment Project Exam Help

## Idea 2 Add WeChat edu\_assist\_pro

Sorted Term	A	C	B
Doc	*	*	*
Cumulative Upper Bound	5	8	10

### Assignment Project Exam Help

if  $\alpha = 9$ , the first document that can score above  $\alpha$  is from \_\_B\_\_

Pivot document is the current d <https://eduassistpro.github.io/>

Is it possible for any doc < PivotDoc to enter int

### Add WeChat edu\_assist\_pro

A	...	2:*	...	UB = 5
B	...	7:*	...	UB = 3
C	...	4:*	...	UB = 2

Pivot term = B

Pivot doc = d7

Case II: smallest DID = PivotDoc

Sorted Term	A	C	B
Doc	7	*	7
Cumulative Upper Bound	5	8	10

C.Doc must be d7

score(d7) could be larger than  $\alpha$

→ fullScore(d7)

→ Adjust  $\alpha$  if necessary

→ all list pointing to d7: next()

# Assignment Project Exam Help

## Idea 2 Add WeChat edu\_assist\_pro

Sorted Term	A	C	B
Doc	*	*	*
Cumulative Upper Bound	5	8	10

A	...	7:*	...	UB = 5
B	...	7:*	...	UB = 3
C	...	7:*	...	UB = 2

if $\alpha = 9$ , the first document that can score above $\alpha$ is from __B__
Pivot document is the current d <a href="https://eduassistpro.github.io/">https://eduassistpro.github.io/</a>

Pivot term = B

Pivot doc = d7

If all term preceding PivotTerm (in the sorted or Add WeChat edu\_assist\_pro n PivotDoc, then PivotDoc is the smallest document that may enter into top-k.

Full scoring:

- Need to “probe” lists that are sorted after the PivotTerm

None in our example, but one can easily add D list.

## Assignment Project Exam Help

WAND Add WeChat edu\_assist\_pro

### 1. Initialization

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



## Assignment Project Exam Help

WAND Add WeChat edu\_assist\_pro

### 2. Finding the Pivot

Assignment Project Exam Help

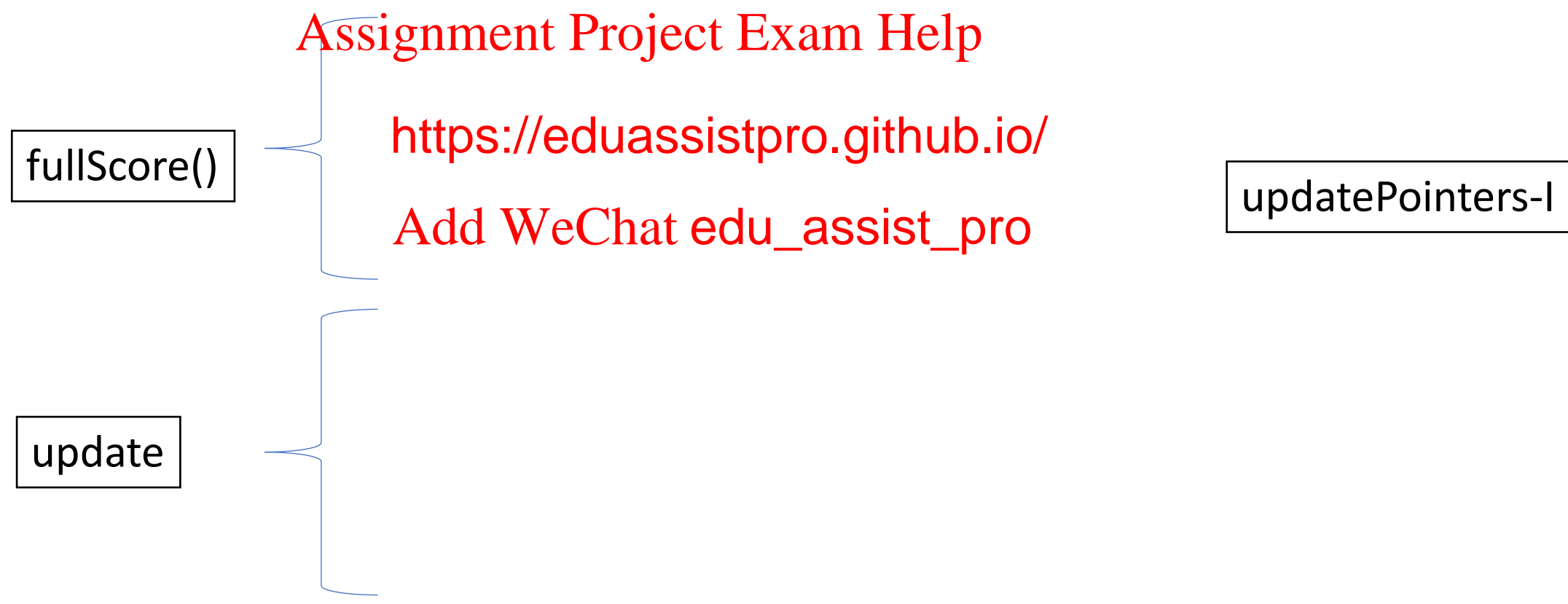
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

WAND Add WeChat edu\_assist\_pro

3a. Case II



## Assignment Project Exam Help

WAND Add WeChat edu\_assist\_pro

3b. Case I

updatePointers-II

It moves all the preceding lists. It is the mWAND optimization for memory-resident index in [Fontoura et al, 2011].

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Example (k=1)

- Step 1:

- $\alpha = 0$

Sorted Term	A	B	C
Doc	1		
Cumulative Upper Bound	4		

- PivotTerm = A
- PivotDoc = 1
- Case II
  - fullScore()  $\rightarrow$  d1 = 11,  $\alpha = 11$
- updatePointers()

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Example (k=1)

- Step 2:

- $\alpha = 11$

Sorted Term	A	B	C
Doc	2		
Cumulative Upper Bound	4		

- PivotTerm = C
- PivotDoc = 2
- Case II
  - fullScore()  $\rightarrow$  d2 = 7,  $\alpha = 11$
- updatePointers()

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Example (k=1)

- Step 3:

- $\alpha = 11$

Sorted Term	C	B	A
Doc	5		
Cumulative Upper Bound	8		

- PivotTerm = B
- PivotDoc = 7
- Case I
  - updatePointers()

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Example (k=1)

- Step 4:

- $\alpha = 11$

Sorted Term	C	B	A
Doc	7		
Cumulative Upper Bound	8		

- PivotTerm = B
- PivotDoc = 7
- Case II
  - fullScore()  $\rightarrow$  d7 = 10,  $\alpha = 11$
- updatePointers()

- Step 5: ...

A	B	C
$UB_A = 4$	$UB_B = 5$	$UB_C = 8$
<1, 3>	<1, 4>	<1, 4>
<2, 4>	<2, 1>	<2, 2>
<7, 1>	<7, 2>	<5, 1>
	<8, 5>	<7, 7>
	<9, 2>	<10, 1>
	<11, 5>	<11, 8>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Assignment Project Exam Help

## Comparisons

	MaxScore	(m)WAND
Pruning Strategy	UB based on fixed ordering of terms → “proactive” pruning	UB based on variable ordering of terms → “passive” pruning
Performance		Better for long queries
Applicability	<a href="https://eduassistpro.github.io/">https://eduassistpro.github.io/</a>	DAT

Add WeChat edu\_assist\_pro



# Assignment Project Exam Help

## Hybrid TAAT DA

- Idea

- Find a good  $\alpha$  (with little cost) and run optimized DAAT algorithm

- $Q = \{A \ B \ C \ D\}$ , in inc length

- $Q1 = \{A, B\}$ ,  $Q2 = Q$

- $\alpha(Q1)$ ,  $L(Q1) = \text{ProcessQuery}(Q1)$

- $L(Q1)$ : documents scored for  $Q1$

- $\text{ProcessQueryDAAT}(Q2; \alpha(Q1), L(Q1))$

- Treat  $L(Q1)$  as another inverted list,  $UB(L(Q1)) = \alpha(Q1)$

# Assignment Project Exam Help

## References

- Efficient Query Evaluation using a Two-Level Retrieval. CIKM 2003.
- Exploring the Magic of WAND. ADCS 2013.
- [Fontoura et al, 2011] or Top-k Queries over Memory-Resident Inv <https://eduassistpro.github.io/>
- Howard R. Turtle, James Flood. Query Strategies and Optimizations. Inf. Process. Manag. 31(6): 831-850 (1995)