# EECS 440 System Design of a Search Engine
## Winter 2021
## Lecture 14: Ranking

Nicole Hamilton
https://web.eecs.umich.edu/~nham/
nham@umich.edu

# Agenda

1. Course details.
2. HW6 and 7 Hashing
3. Ranking.

# Agenda

1. Course details.
2. HW6 and 7 Hashing
3. Ranking.

# details

1.  A few teams have made good progress (one finished) on the hashtable and top 10.  No one has started the hashblob and hashfile.

2.  Reading:

    a.  My paper on dynamic ranking.

    b.  Brian Fung, *"Here's what we know about Google's mysterious search engine"*, The Washington Post, August 28, 2018.

    c.  Look at (no need to read in any detail) Google's *Page Quality Rating Guidelines*.

    d.  Marc Najork and Allan Heydon, "*Mercator*", September 26, 2001.

# Agenda

1. Course details.

2. HW6 and 7 Hashing

3. Ranking.

In Homework 6, you will build a conventional hash table.  If you build it with -DVerbose, you get timing information. I also give you some sample input.

```
tcsh-31% make verbose
g++ -DVerbose Top10.cpp Common.cpp TopN.cpp -o Top10
g++ -DVerbose HashTable.cpp Common.cpp -o HashTable
g++ -DVerbose HashBlob.cpp Common.cpp -o HashBlob
g++ -DVerbose HashFile.cpp Common.cpp -o HashFile
tcsh-32% wc BigJunkHtml.txt
  62001  215994 2209965 BigJunkHtml.txt
tcsh-33%
```

In Homework 6, you will build a conventional hash table.  If you build it with -DVerbose, you get timing information. I also give you some sample input.

```
tcsh-33% ./HashTable BigJunkHtml.txt
Number of tokens = 215994
Total characters = 1782086
Average token length = 8.25063 characters

Building HashTable
Elapsed time = 14759200 ticks

Optimizing HashTable
Elapsed time = 591300 ticks

Enter search words:
hello world how are you
88   hello
43   world
91   how
650   are
675   you
Elapsed time = 7696912200 ticks

tcsh-34%
```

Here's the top 10.

```
tcsh-34% ./Top10 BigJunkHtml.txt
Number of tokens = 215994
Total characters = 1782086
Average token length = 8.25063 characters

11931   <li><a
6605    the
3314    <a
3144    to
3088    a
2223    and
2059    of
1930    C
1837    is
1768    </td>
tcsh-36%
```

In HW7, you will build a HashBlob in memory and then search it.

```
tcsh-35% ./HashBlob BigJunkHtml.txt
Number of tokens = 215994
Total characters = 1782086
Average token length = 8.25063 characters

Building HashTable
Elapsed time = 14677300 ticks

Optimizing HashTable
Elapsed time = 575700 ticks

Building HashBlob
Elapsed time = 2503900 ticks

HashBlob size = 942840 bytes

Enter search words:
hello world how are you
88    hello
43    world
91    how
650    are
675    you
Elapsed time = 6315604700 ticks
```

You will also build a HashBlob in as a mapped file.

```
tcsh-36% ./HashBlob BigJunkHtml.txt Blob.bin < /dev/null
Number of tokens = 215994
Total characters = 1782086
Average token length = 8.25063 characters

Building HashTable
Elapsed time = 14835000 ticks

Optimizing HashTable
Elapsed time = 637600 ticks

Building HashFile = Blob.bin
Elapsed time = 18153000 ticks

HashBlob size = 942840 bytes
Elapsed time = 24100 ticks

tcsh-37%
```

You will then search the HashBlob in as a mapped file. (The elapsed time reflects that I typed the input search words!)

```
tcsh-37% ./HashFile Blob.bin
Loading HashBlob from Blob.bin
Elapsed time = 105200 ticks

HashBlob size = 942840 bytes

Enter search words:
hello world how are you
88    hello
43    world
91    how
650    are
675    you
Elapsed time = 72123829500 ticks

tcsh-38%
```

# Agenda

1. Course details.
2. HW6 and 7 Hashing
3. Ranking.

# Ranking

Objective is to order pages the same way a human would do using software.

Do a calculation based on what's known about:

1. All the documents and words in the index.

2. That page.

3. Match between the query and the page.

# Ranking

1. The rank value should obey the desired ordering relationship, that a better page will get a better score.

2. Since search engines typically broadcast a query to a large number of machines with differing fractions of the web and then combine the results, the calculation should be consistent even though the samples might be a little different.

3. Beyond ordering, the value is otherwise meaningless: If page A's rank is twice that of page B, it does not mean that A is twice as good as B.

Only two ways to get better at ranking:

1. Have more or better information.
2. Make better use of it.

# Static vs. dynamic rank

*Static*        Quality of the page independent of the query, e.g., PageRank, length of the URL, title or page, domain (.gov or .edu vs. .biz), whether it contains images, pornographic content, etc.

*Dynamic*    Quality of a page as possible result for a specific query considering both static rank and the quality of the match between the query and the page.

# Static rank

Some pages are just better than others before you know anything about the query.

# Static rank

Some domains are better than others, e.g., .gov or .edu over .biz.

Short URLs are better.

Short titles are probably better.

Some pages may be obvious spam.

Some pages may have lots of other pages pointing to them.

# PageRank

A detour into the world's most famous link-analysis algorithm.

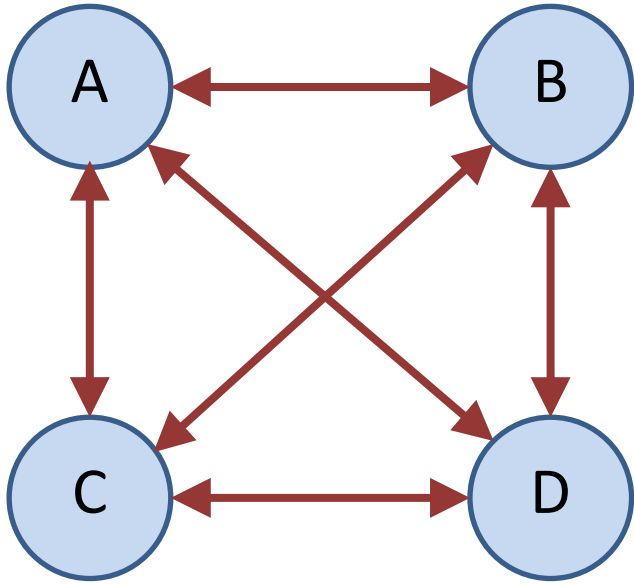The basic idea: The more and better links to a page, the more likely it should rank higher.

# PageRank random surfer

The model was that people surfed the web, somewhat randomly either clicking on one of the links on the page or going somewhere else.

If an important page pointed to yours, some of that importance should bleed onto yours.

# Basic PR algorithm

1. PR output is a vector of probabilities that a person randomly clicking links will arrive at a particular page.

2. Initially all probabilities usually assumed equal (or maybe not!)

3. Links from a page to itself are ignored.

4. Multiple links from one page to another are treated as a single link.

5. The PR transferred from one page to another is its PR divided by number of pages it links to.

6. At each iteration, the new PR of a given page is calculated as the sum of the PRs transferred to it.

7. Repeat until it settles.

$$PR(A) = \frac{PR(B)}{3} + \frac{PR(C)}{3} + \frac{PR(D)}{3}$$

$$PR(B) = \frac{PR(A)}{3} + \frac{PR(C)}{3} + \frac{PR(D)}{3}$$

$$PR(C) = \frac{PR(A)}{3} + \frac{PR(B)}{3} + \frac{PR(D)}{3}$$

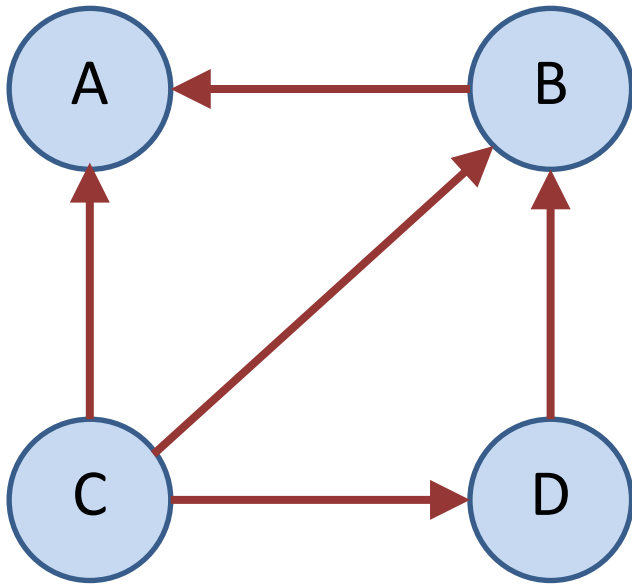$$PR(D) = \frac{PR(A)}{3} + \frac{PR(B)}{3} + \frac{PR(C)}{3}$$

$$PR(A) = \frac{PR(B)}{1} + \frac{PR(C)}{3}$$

$$PR(B) = \frac{PR(C)}{3} + \frac{PR(D)}{1}$$

$$PR(C) = 0$$

$$PR(D) = \frac{PR(C)}{1}$$

$$PR(u) = \sum_{u \in B_u} \frac{PR(v)}{L(v)}$$
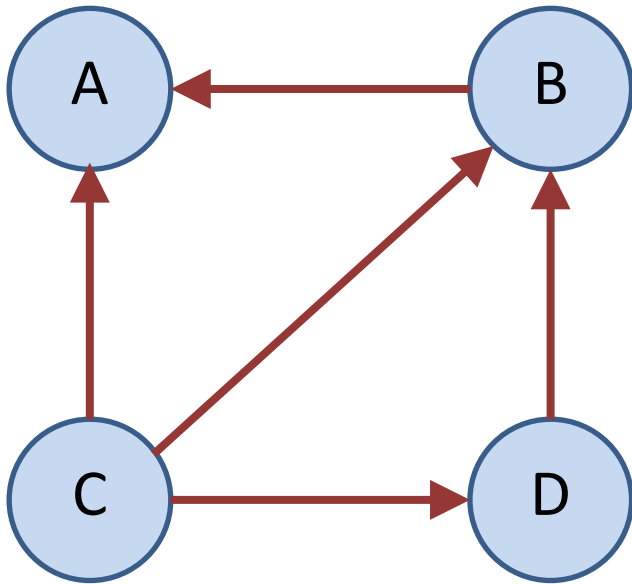
Where
  $u$ is a page,
  $PR(u)$is is the PageRank of page $u$,
  $B_u$ is the set of all pages that link to $u$,
  $L(v)$ is the number pages linked from $v$.

# Damping factor

But in this basic version, PR sinks could happen, where at every iteration, a site just got more and rank.

They solved this by adding the notion that their imaginary surfer randomly clicking links will eventually stop clicking and simply start over at some other random page.  This became a damping factor in PageRank.

$$PR(u) = \frac{1-d}{N} + d \sum_{v \in B_u} \frac{PR(v)}{L(v)}$$

Where
    $u$ is a page,
    $PR(u)$is is the PageRank of page $u$,
    $B_u$ is the set of all pages that link to $u$,
    $L(v)$ is the number pages linked from $v$,
    $d$ is the damping factor, typically ~0.85.

# PageRank

It obviously did work and Google got better results.

It also gave halo of special legitimacy to their results, that they were scientific and unbiased.

# At Microsoft

We gamely expected our version of PageRank to represent about half the overall rank value, largely based on the hype around it.

Turned out it was very expensive to calculate and represented only a small part of the final rank score.

Mark Najork of [Mercator](#) fame argued for lumping whole domains together in something called DomainRank.  But it hadn't yet worked when I left.

# Ranking process

1. Compile the query.

2. Search the index for matching pages.

3. Return a list of the n best with scores indicating estimated quality.

4. May also return debug information to allow the scoring calculation to be examined.

# Question

1. How should you find the n best?

2. Should you get the entire list and then sort?

# Question

1. How should you find the n best?
2. Should you get the entire list and then sort?

No, you should probably insertion sort into array of n elements.

# Simple search engine query language

```
<Constraint>          ::= <BaseConstraint> { <OrOp> <BaseConstraint> }

<OrOp>                ::= 'OR' | '|' | '||'

<BaseConstraint>      ::= <SimpleConstaint> { [ <AndOp> ] <SimpleConstraint> }

<AndOp>               ::= 'AND' | '&' | '&&'

<SimpleConstraint>    ::= <Phrase> | <NestedConstraint> |
                          <UnaryOp> <SimpleConstraint> |
                          <SearchWord>

<UnaryOp>             ::= '+' | '-' | 'NOT'

<Phrase>              ::= '"' { <SearchWord> } '"'

<NestedConstraint>   ::= '(' <Constraint> ')'
```
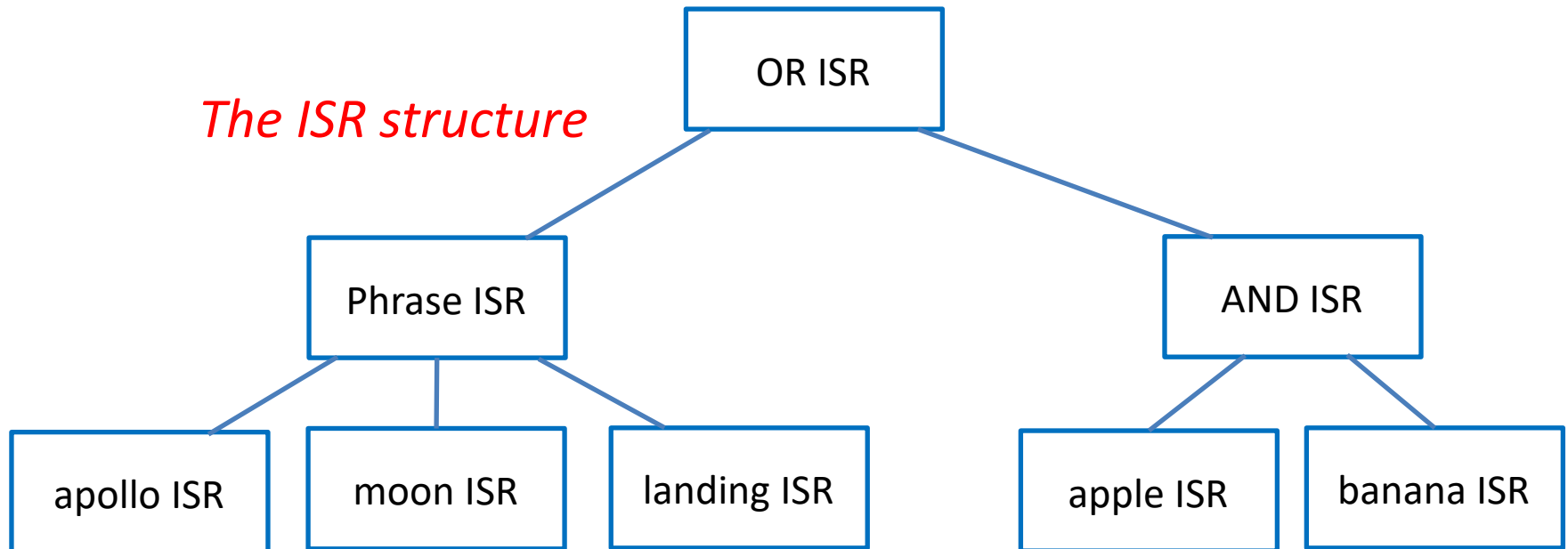
# The query language and the ISRs can be recursive

"apollo moon landing" | ( apple banana )

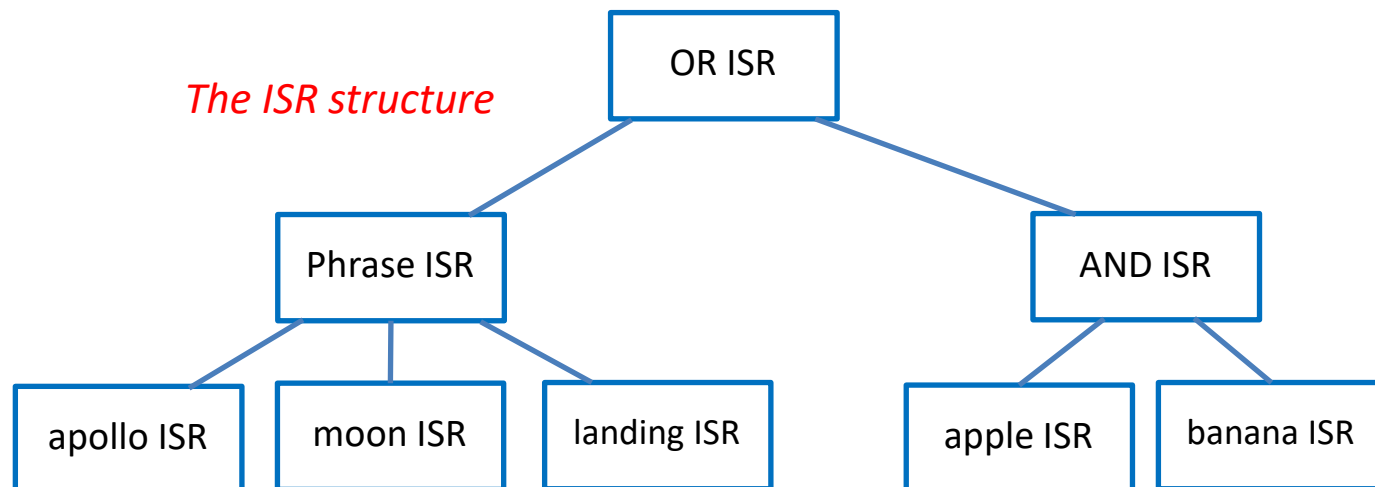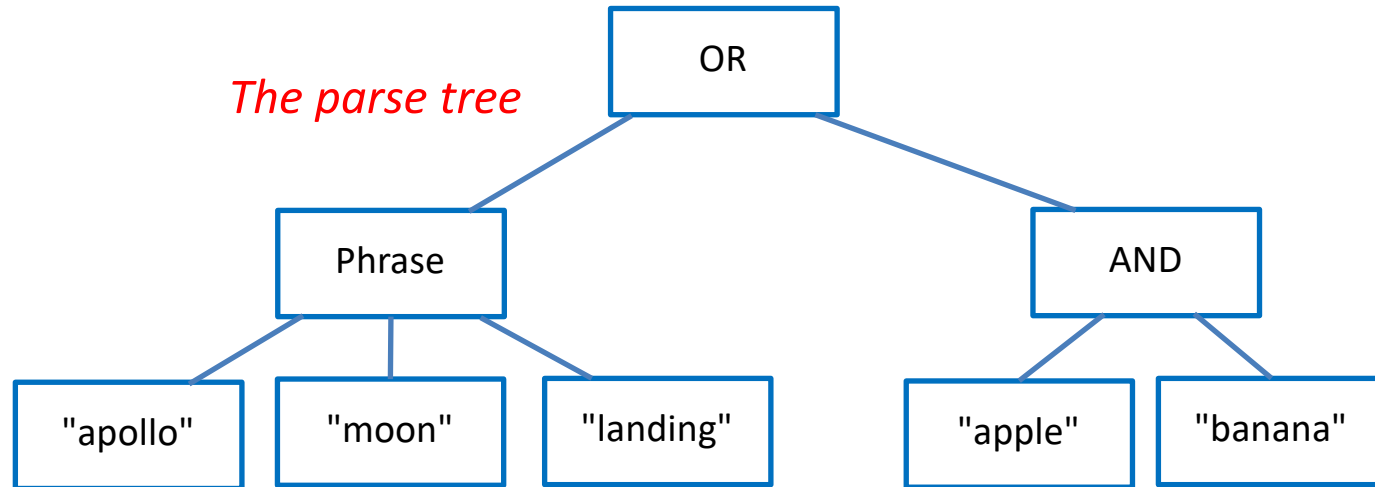*The parse tree*

```
                              OR
                     /                  \
                Phrase                   AND
              /    |    \               /     \
      "apollo" "moon" "landing"   "apple"   "banana"
```
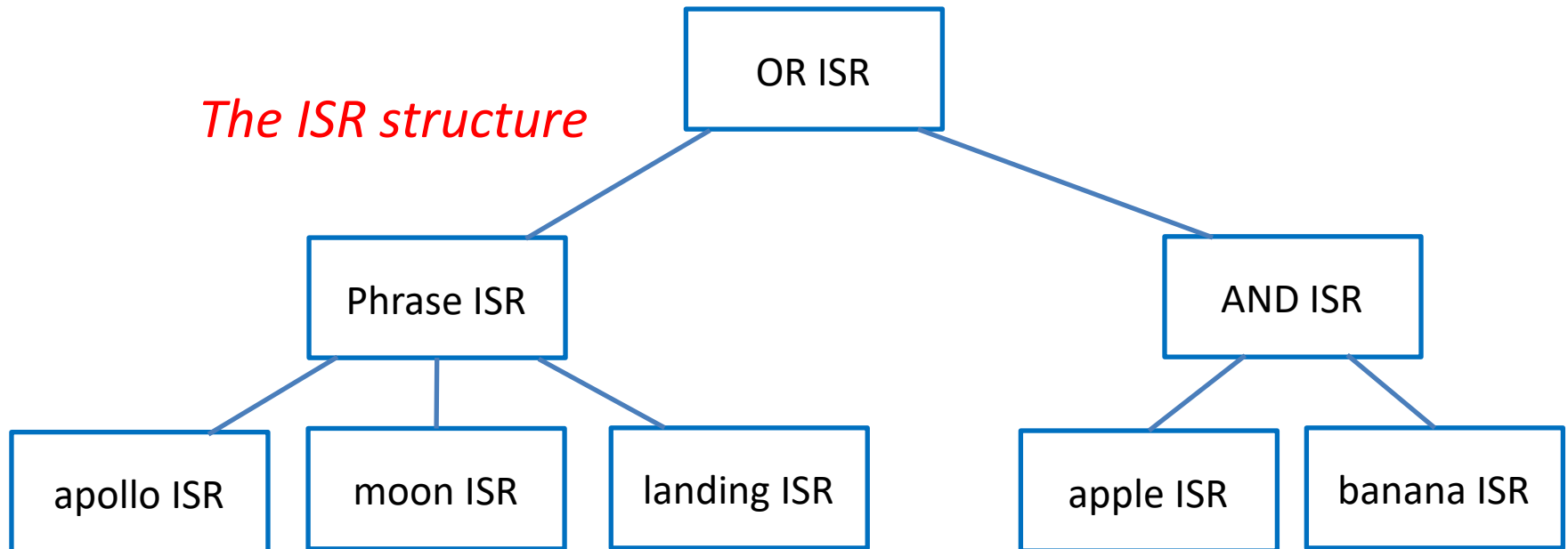
# The query language and the ISRs can be recursive

"apollo moon landing" | ( apple banana )

*The ISR structure*

"apollo moon landing" | ( apple banana )

The trees are the same.

*The parse tree*

```
                          OR
             ┌────────────┴────────────┐
          Phrase                      AND
      ┌─────┼─────┐                  ┌──┴──┐
  "apollo" "moon" "landing"      "apple" "banana"
```

*The ISR structure*

```
                        OR ISR
             ┌────────────┴────────────┐
        Phrase ISR                  AND ISR
      ┌─────┼─────┐                  ┌──┴──┐
 apollo ISR moon ISR landing ISR  apple ISR banana ISR
```

Technically, is this correct?

"apollo moon landing" | ( apple banana )

*The ISR structure*

```
                           ┌──────────┐
                           │  OR ISR  │
                           └──────────┘
                          /            \
                 ┌───────────┐        ┌──────────┐
                 │ Phrase ISR│        │ AND ISR  │
                 └───────────┘        └──────────┘
                 /    |    \           /        \
        ┌─────────┐ ┌────────┐ ┌──────────┐  ┌──────────┐ ┌───────────┐
        │apollo ISR│ │moon ISR│ │landing ISR│ │apple ISR │ │banana ISR │
        └─────────┘ └────────┘ └──────────┘  └──────────┘ └───────────┘
```

*How many ISRs does it really take to do this?*

# Decorating

Add characters that get stripped out during HTML parsing to indicate special characteristics or types of posts, e.g.,

| | |
|---|---|
| amazon | amazon in the body text |
| #amazon | amazon only in the URL |
| @amazon | amazon only in the title |
| $amazon | amazon only in the anchor text |
| % | End-of-document token. |

Might also be used for *stemming*:

| | |
|---|---|
| swim* | swim, swims, swimming, etc. |

# Decorating vs. attributes

Use decorating when you'd like to use it for searching, to shrink the size of a post or because you'd like to separate the scoring for hits in the title vs. the body for example.

Use attributes when the ranker will want the information about each post and it could be different every time.

The AND of apple and banana might actually take 12 ISRs if the terms are decorated, e.g., @ = anchor, # = title, / = url.

"apollo moon landing" phrase actually requires 4 ISRs for each stream (anchor, title, URL and body) + an OR ISR = 17 ISRs.

Assume decorations:  @ = anchor, # = title, / = url

# Coming back to the question, is this correct?

"apollo moon landing" | ( apple banana )



*1 ISR* — OR ISR

*17 ISRs* — Phrase ISR

*12 ISRs* — AND ISR

apollo ISR     moon ISR     landing ISR

apple ISR     banana ISR

*If the terms are decorated, it could take 30 ISRs.*

Matches in order of importance:

1. Anchor text.
2. URL.
3. Title.
4. Body.

Exception is Japan were URL matches are less useful due to mismatch between Kanji or Katakana queries and transliterated URLs.

URLs may need dictionary word-breaking or regex-style matching to be useful.  Again, beyond the scope.

Does it matter how many ISRs you use?

How much effect will this have on query search time?

What dominates search time?

# Basic ranking

1.  Matching pages found by the constraint solver but that only finds the page and the static information about the page from the enddoc.

2.  Queries are flattened.  (Very hard to estimate the probability of phrases or other combinations of OR'ing and AND'ing terms.)

3.  ISRs are reset to the beginning of the document, then advanced through the page, extracting data about where the search words were found.

4.  Three strategies from there.

# Queries are flattened

These queries all match different sets of pages:

apollo moon landing

( apollo | moon ) landing

"apollo moon landing"

But for scoring, they're all flattened to the same list of search words.

# Three strategies

1.  Bag of words:  The more hits the better.

2.  Heuristics:  Look for exact matches, matches in the right places, hand or machine-tuned.

3.  Machine learning, typically with a neural net.

# Bag of words

Count the number of matches of each of the search words, typically weighted by the frequency of the word within the corpus.

Two most famous:

1. tf-idf
2. BM25

# tf-idf

Term-frequency, inverse document frequency.

Bag of words technique. The more occurrences of a rare word, the better.

Combined:

1. Term weighting based on frequency invented by to Hans Peter Luhn in 1957.

2. Statistical interpretation invented by Karen Spärck Jones in 1972.

Term frequency

$$TF(t) = \frac{Number\ of\ times\ t\ appears\ in\ a\ document}{Total\ number\ of\ terms\ in\ the\ document}$$

Inverse document frequency

$$IDF(t) = \frac{\log_e(Total\ number\ of\ documents)}{Number\ of\ documents\ with\ term\ t}$$

Tf-idf is the product. The more the better.

$$Tf - idf(t) = TF(t) * IDF(t)$$

# Okapi BM25

Okapi Best Matching function, a similar bag of words technique.

$$score(D, Q) = \sum_{i=1}^{n} IDF(q) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgdl})}$$

$f(q_i, D)$ is $q_i$'s term frequency in Document D.

$|D|$ is the length of the document in words.

$avgdl$ is the average document length.

$k_1$ and $b$ are free parameters you get to choose. Typically

$$k_1 \in [1.2, 2.0]$$

$b$ = 0.75

# Bag of words

Problem is, bag of words techniques simply don't work very well on the web because they don't do well at distinguishing quality, especially, to find the best match.

They simply cannot distinguish that a page with all the search words in the right order, as an exact phrase, or near the top of the page is better than a page where the words are randomly scattered.

# Bag of words

Here's a sample NY Times page from Jan 25, 2021.

# Bag of words

Here it is stripped of HTML and CSS but the text remains.

# Bag of words

Here it is with the words in the title and body sorted.  Tf-idf can't tell the difference.

# Heuristics

1. Incrementing a flurry of low-level counters, e.g., number of times each word occurred, number of times an exact phrase matching the query was found.

2. Large number of cheap heuristics that are expected to provide evidence of the quality of the match.

# Query processing

1. Compile the query into a structure of ISRs.

2. Pass the ISRs to the constraint solver to find matching pages.

3. For each matching page:

   a. Move the ISRs back to the beginning of the page and scan for hits.

   b. Calculate a rank value and insertion sort the page into a list of n best.

Example phrase query "quick brown fox" and the their 3 posting lists.

Individual search words appear individually many times but there are only two exact phrases. The constraint solver will stop on the first phrase match, then it's up to the ranker to decide what next.

|  | brown |  |
| --- | --- | --- |
|  | 83 |  |
|  | 94 |  |
|  | 170 |  |
| quick | 179 |  |
| 62 | 216 |  |
| 69 | 227 |  |
| 84 | 400 | fox |
| 311 | 417 | 284 |
| 421 | 422 | 423 |
| 430 | 516 | 580 |
| 559 | 795 | 612 |
| 619 | 826 | 796 |
| 794 | 828 | 912 |
| 952 | 957 | 958 |

|  | brown |  |
| --- | --- | --- |
|  | 83 |  |
| quick | 94 |  |
| 62 | 170 |  |
| 69 | 179 |  |
| 84 | 216 |  |
| 311 | 227 | fox |
| 421 | 400 | 284 |
| 430 | 417 | 423 |
| 559 | 422 | 580 |
| 619 | 516 | 612 |
| 794 | 795 | 796 |
| 952 | 826 | 912 |
|  | 828 | 958 |
|  | 957 |  |

Example phrase query "quick brown fox".

A suggested first step in ranking is to flatten the query, pulling out all the individual ISRs and seeking them to the beginning of the document.

| quick | brown | fox |
|-------|-------|-----|
|       | 83    |     |
|       | 94    |     |
|       | 170   |     |
|       | 179   |     |
| 62    | 216   |     |
| 69    | 227   |     |
| 84    | 400   |     |
| 311   | 417   | 284 |
| 421   | 422   | 423 |
| 430   | 516   | 580 |
| 559   | 795   | 612 |
| 619   | 826   | 796 |
| 794   | 828   | 912 |
| 952   | 957   | 958 |

| quick | brown | fox |
|-------|-------|-----|
| 62    | 83    | 284 |
| 69    | 94    | 423 |
| 84    | 170   | 580 |
| 311   | 179   | 612 |
| 421   | 216   | 796 |
| 430   | 227   | 912 |
| 559   | 400   | 958 |
| 619   | 417   |     |
| 794   | 422   |     |
| 952   | 516   |     |
|       | 795   |     |
|       | 826   |     |
|       | 828   |     |
|       | 957   |     |

Example phrase query "quick brown fox".

From there, you can move the ISRs any way you like to extract data as long as they only go forward.

Word counts:

| | | quick | brown | fox |
|---|---|---|---|---|
| quick | 10 | 62 | 83 | 284 |
| brown | 14 | 69 | 94 | 423 |
| fox | 7 | 84 | 170 | 580 |
| | | 311 | 179 | 612 |
| | | 421 | 216 | 796 |

Two possible strategies:
1. Simply count the words.
2. Look for places where the words occur together.
3. Possible combinations of the three words = 10 * 14 * 7 = 980 for a very short document and will grow with longer queries.
4. Not possible to visit all combinations if all the ISRs only go forward.
5. May either read the lists in or process on the fly.

| quick | brown | fox |
|---|---|---|
| 430 | 227 | 912 |
| 559 | 400 | 958 |
| 619 | 417 | |
| 794 | 422 | |
| 952 | 516 | |
| | 795 | |
| | 826 | |
| | 828 | |
| | 957 | |

Example phrase query "quick brown fox".

I'm going to show you a simpler strategy I used at Microsoft.



| quick | brown | fox |
| --- | --- | --- |
| 62 | 83 | 284 |
| 69 | 94 | 423 |
| 84 | 170 | 580 |
| 311 | 179 | 612 |
| 421 | 216 | 796 |
| 430 | 227 | 912 |
| 559 | 400 | 958 |
| 619 | 417 | |
| 794 | 422 | |
| 952 | 516 | |
| | 795 | |
| | 826 | |
| | 828 | |
| | 957 | |

(57) **ABSTRACT**

A system and a method for locating and presenting electronic documents most-likely of interest to the user. A plurality of search terms to be located in a set of electronic documents is received. One of the search terms is selected as the anchor term, and occurrences of the anchor term are located within the documents. For each located occurrence of the anchor term, a set of search term occurrences is selected. These sets include an occurrence of each search term, and the occurrences are selected by choosing the search term occurrences that are closest to a desired placement for the search terms. With each set of search terms, the method associates a value indicating the extent to which the selected occurrences vary from the desired placement. The electronic documents are ranked and presented to the user in accordance with this value. The invention further includes systems and methods for locating and presenting Web pages and for searching the Internet.

40 Claims, 5 Drawing Sheets

I was very concerned that my ranker would single-handedly blow us out of the water on perf and I wanted something really simple.

It was all integer math and done in a single pass over each page.

# AND'ing streams

| quick | 10 27 105 | 513 518 520 | |
| brown | 28 50 62 70 | 514 | 790 |
| fox | 87 106 | 515 550 | 1200 |
| #DocEnd | 112 | 570 1006 | 1704 |

quick fox     *How many possible combinations?  6*
*Can you reach all of them in a single pass, all ISRs only moving*
*forward?  No.*

Would prefer a technique that allows the ISRs to be moved only in a single pass.

# Basic strategy

Choose rarest word as an anchor and advance that ISR through each occurrence in the document.

At each occurrence of the rarest word, advance the other ISRs to position them as close as possible to the desired position in the flattened list.

Requires only one stage of look-ahead.

# Spans

1.  If there is at least one occurrence of each query term, that's a "span" which can be further distinguished as ordered, short, an exact phrase, etc.

2.  Only the rarest word is guaranteed to a unique occurrence.

3.  The other words may be reused.

4.  Also count doubles and triples, meaning various combinations of just 2 or 3 of the words in the query, one of which must be the rarest.

5.  Parameterized threshold for short vs. long spans, frequent vs. infrequent, etc.

6.  Most features are binary and either occur or do not occur unambiguously.

Example phrase query "quick brown fox".

What I did was *pick the rarest word* and *then arrange the other ISRs* to as close as possible to their desired locations in the query relative to each occurrence of the rarest word.

| quick | brown | fox |
|-------|-------|-----|
| 62    | 83    | 284 |
| 69    | 94    | 423 |
| 84    | 170   | 580 |
| 311   | 179   | 612 |
| 421   | 216   | 796 |
| 430   | 227   | 912 |
| 559   | 400   | 958 |
| 619   | 417   |     |
| 794   | 422   |     |
| 952   | 516   |     |
|       | 795   |     |
|       | 826   |     |
|       | 828   |     |
|       | 957   |     |

Word counts:
quick           10
brown           14
fox              7

This reduces the number of combinations to be scored in this example from 10*14*7 = 957 to 7 and can be done a single pass with a single stage of lookahead.

Call each combination a *span*.

Example phrase query "quick brown fox".

Iterate over the rarest word hits, arranging the other ISRs to as close as they could be to their desired locations in the query. Here are the first 3 *spans*.

Example phrase query "quick brown fox".

For each combination, decide if it's an *exact phrase* or all the words *in order*, or *close together*, incrementing an associated counter. Shown here, there is one exact phrase but none of the rest are in order or particularly close together.

Example phrase query "quick brown fox".

Example phrase query "quick brown fox".

This last one is *close together* (say, less than 10 words apart) and *in order* but it's *not an exact phrase*.

|  | brown |  |
|---|---|---|
|  | 83 |  |
|  | 94 |  |
|  | 170 |  |
| quick | 179 |  |
| 62 | 216 |  |
| 69 | 227 |  |
| 84 | 400 |  |
| 311 | 417 | fox |
| 421 | 422 | 284 |
| 430 | 516 | 423 |
| 559 | 795 | 580 |
| 619 | 826 | 612 |
| 794 | 828 | 796 |
| 952 | 957 | 912 |
|  |  | 958 |

Example phrase query "quick brown fox".

Set some thresholds and accumulate some counts,
which can scored at the end.

brown

| | |
|---|---|
| | 83 |
| | 94 |
| | 170 |
| quick | 179 |

| Thresholds | Values |
|---|---|
| Max to be short | 10 |
| Min to be frequent | ? |
| Min to be most | ? |
| Min to be near the top | ? |

| quick | brown |
|---|---|
| 62 | 216 |
| 69 | 227 |
| 84 | 400 |
| 311 | 417 |
| 421 | 422 |
| 430 | 516 |
| 559 | 795 |
| 619 | 826 |
| 794 | 828 |
| 952 | 957 |

fox

| fox |
|---|
| 284 |
| 423 |
| 580 |
| 612 |
| 796 |
| 912 |
| 958 |

| Heuristic | Count | Weight | Score |
|---|---|---|---|
| Number of short spans | 1 | 5 | 5 |
| Number of in order spans | 1 | 2 | 2 |
| Exact phrases | 2 | 10 | 20 |
| Number of spans near the top | ? | ? | ? |
| All word are frequent | ? | ? | ? |
| Most words are frequent | ? | ? | ? |
| Some words are frequent | ? | ? | ? |

| *Total dynamic rank score* | ? |
|---|---|

Example phrase query "quick brown fox".

One reason to decorate words as *anchor*, *URL*, *title* or *body* is so they can be separated in to separate *streams* (separate sets of ISRs) and *scored separately* with the same algorithm but different scoring parameters.

| Stream | Weight | Score |
|---|---|---|
| Anchor | ? | ? |
| URL | ? | ? |
| Title | ? | ? |
| Body | ? | ? |
| | | |
| *Total dynamic rank score* | | ? |

Each of the streams can start at location 1 relative to the start of the document. The document length is the length of the longest stream.

Not all streams may contain all the words.

brown

83
94
170
179

| quick | brown |
|---|---|
| 62 | 216 |
| 69 | 227 |
| 84 | 400 |
| 311 | 417 |
| 421 | 422 |
| 430 | 516 |
| 559 | 795 |
| 619 | 826 |
| 794 | 828 |
| 952 | 957 |

fox

284
423
580
612
796
912

958

Example phrase query "quick brown fox".

We do a similar heuristic calculation of the *static rank*,
the quality of the page *independent of the quer*y.

| Heuristic | Weight | Score |
|---|---|---|
| Short title | ? | ? |
| Nice document length | ? | ? |
| Short URL | ? | ? |
| Lots of anchor text | ? | ? |
| edu/gov/com/etc domain | ? | ? |
| PageRank if known | ? | ? |
| | | |
| Total static rank score | | ? |

brown
83
94
170
179

quick | 216
62 | 227
69 |
84 | 400

fox
311 | 417 | 284
421 | 422 | 423
430 | 516 | 580
559 | 795 | 612
619 | 826 | 796
794 | 828 | 912
952 | 957 | 958

Example phrase query "quick brown fox".

Combine static and dynamic rank to get a final score.

| Component | Weight | Score |
|-----------|--------|-------|
| Static rank | ? | ? |
| Dynamic rank | ? | ? |
| | | |
| *Total rank score* | | *?* |

The result is a *linear combination* of features. We're just adding them up.

You should be able to achieve reasonably good results with a reasonably small number of heuristics and simple *hand-tuning* starting from some rough guesses at, e.g., the relative importance of an exact phrase versus lots of individual hits.

brown

83
94
170

quick        179

62        216
69        227
84        400

                                    fox
311        417
                                    284
421        422
                                    423
430        516
                                    580
559        795
                                    612
619        826
                                    796
794        828        912

952        957        958

The actual score is calculated as a linear combination of features, which may be thought of as:

$$R = \Sigma\ C_i(Q) * A_i(P, F_i) * S_i(F_i)$$

Where:

$R$ = Overall Rank

$Q$ = The Query and its characteristics, e.g., the number of rare vs. common words in the query.

$P$ = The Page and its characteristics, e.g., the number of words in the URL or title.

$F_i$ = An arbitrary feature observation, e.g., an exact phrase in the title, or a raw value, e.g., raw PageRank.

$S_i$ = Scaling for feature $F_i$ from the raw number space of the feature into a nominal 0.0 .. 1.0 range.

$A_i$ = Attribute scaling for feature $F_i$, possibly dependent on the characteristics of the page $P$.

$C_i$ = Coefficient for feature $F_i$, depending on the characteristics of the query $Q$.

$F_i$ = $E_i(Q, P, T)$

# Steps to ranking

1. Decide what information to collect and how to measure it.

2. Decide how to measure the quality of the result.

3. Pick a method for scoring a page based on the inputs.

4. Tune the system by testing it on sample queries and adjusting parameters, collecting more information or changing the scoring algorithm.

# Measuring quality

At Microsoft, initially, I just ran queries and eyeballed the results and fiddled with the parameters.

Results were surprisingly good.

# Measuring quality

We were also collecting labeled pages.  We scraped queries and results from several engines and then paid people to rate results on a 0 to 5 scale.

5    Definitive result.  Unlikely any other page could be a better result for this query, e.g., whitehouse.gov for "whitehouse".
:
0    Completely deleterious result, e.g., porno, spam, phishing.

This is similar to what Google is doing with their Page Quality Rating Guidelines, creating specificity in the ratings.

Altogether, I think we had about 50K queries and about 150K labeled results.

# Measuring quality

Next came a tool allowing us to do side-by-side comparison of results with current best parameters in a frame on the left and with a set you could tweak on the right.

If the pages had been labeled, it reported the quality and attempted to score the resulting set.

We ran a competition to see who could come up with the best settings.

# Measuring quality

Finally, we added a tuner that could adjust the parameters mechanically by gradient descent:  Tweak an individual parameter rerun all the searches and measure whether the results got better (more highly ranked pages for each query.)

But we had lots of problems in the methodology of what to do with unlabeled pages.

Initially they were assumed to be "average".   Later, a "promising proximity" heuristic was added to the tuner bump the estimate for unlabed pages if the search words were found close together.

The effect was to tune my complex ranker ended up being tuned to behave like the tuner's naïve ranker.

Better strategy would have been to only tune based on labeled results, discarding any unlabeled results.

# Generating labeled results

You can't afford to pay people to label results  but you might assume that Google is pretty good, so you might simply try to get pages in the same order as Google.
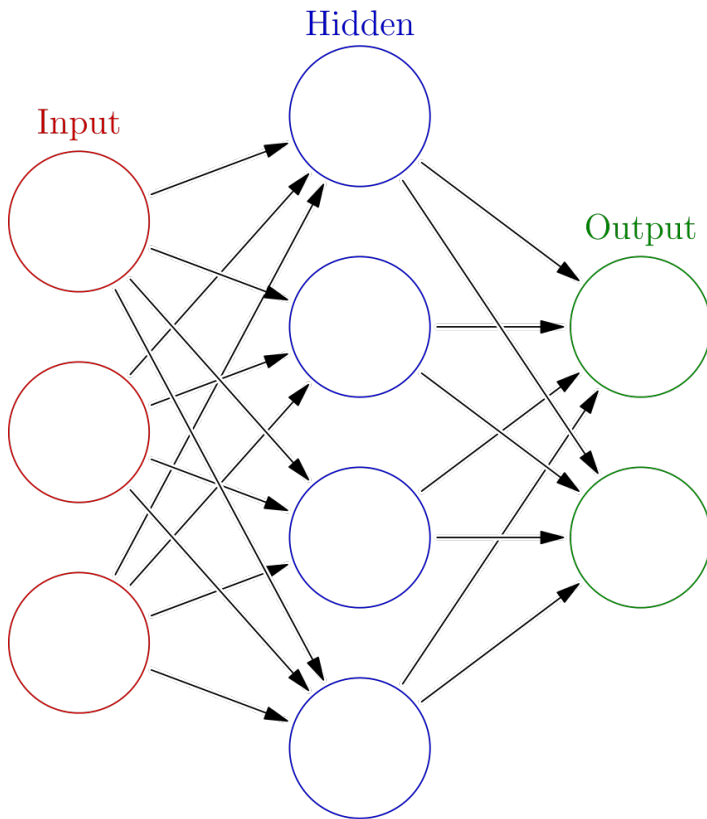
Your ranker won't be as good as Google's, so it might not be able to discern the difference between the top results on the first page, but perhaps it could distinguish between the first result on the first page and the first result on the fifth page.

# How to do better

To do better than heuristics, you will probably need a neural network.

Example phrase query "quick brown fox".

A machine learning strategy would be to simply collect all the lists and pass everything to a neural network, which must be trained. Beyond the scope here.



| quick | brown | fox |
|-------|-------|-----|
| 62 | 83 | 284 |
| 69 | 94 | 423 |
| 84 | 170 | 580 |
| 311 | 179 | 612 |
| 421 | 216 | 796 |
| 430 | 227 | 912 |
| 559 | 400 | 958 |
| 619 | 417 | |
| 794 | 422 | |
| 952 | 516 | |
| | 795 | |
| | 826 | |
| | 828 | |
| | 957 | |