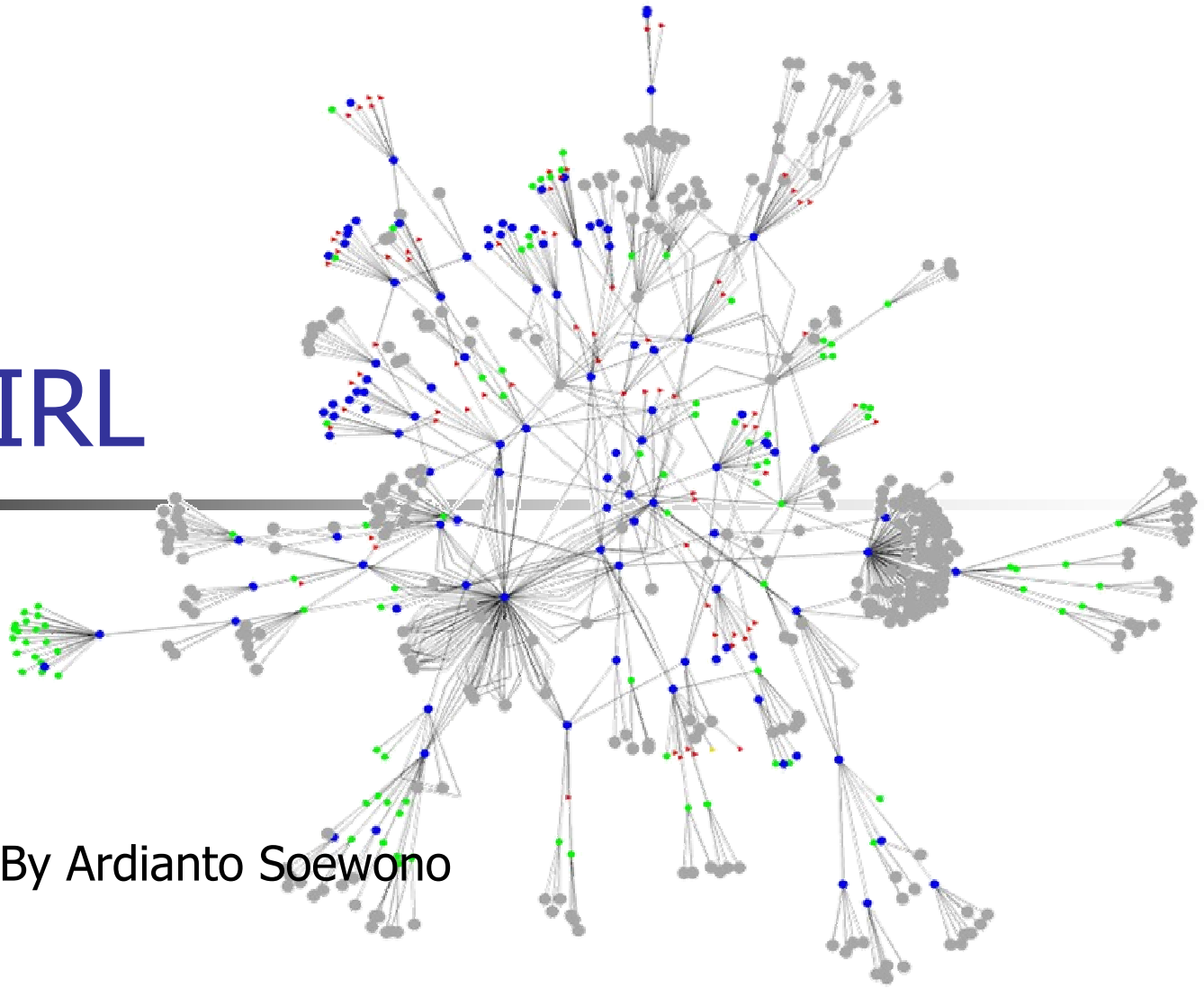




WHIRL

Presented By Ardianto Soewono





What's WHIRL?

- WHIRL stands for Word-based Heterogeneous Information Retrieval Logic.
- WHIRL is a probability database logics used to integrate database that doesn't have common domains.
- As reflected on the name, WHIRL uses words-matching procedure to join tables.



Why WHIRL?

- There is no common domain in the heterogeneous databases.
Exp : “Microsoft” and “Microsoft Corp”, “AT&T” and “AT&T Lab”
- Mapping from local domain to global domain is hard.
The mapping from name constants to real entity can be different from databases to databases
- High accuracy is needed in order to avoid erroneous or missing answers to user queries.



Document Vector

- WHIRL uses document vector as data model
 - There is set of vocabulary terms T .
 - A fragment of text will be represented by a document vector, each component corresponds to a term $t \in T$
 - Each component of document vector v which corresponds to a term t will be denoted as v^t
 - The weight scheme is as follow :
 1. If v^t is not present in the document represented by v , weight = 0
 2. Else :

$$V^t = (\log(TF_{v_t}) + 1). \log(IDF_t)$$

Where TF_{v_t} = the number of items term t occurs in document v

$$IDF_t = |C| / |C_t|$$

where C_t = subset of documents C that contains term t



Document Vector (cont.)

- The similarity of 2 document vectors v and w can be computed as follows :

$$sim(v, w) = \sum_{t \in T} v^t \cdot w^t$$

- The idea of this scheme is giving more weight to infrequent terms.

Exp : Terms "Inc", "Ltd" will have less weight than "AT&T" and "Microsoft".



Scoring

- An EDB consists of a term vocabulary T and set of relations $\{p_1, p_2, \dots, p_n\}$.
- Each relation p with arity k has set of tuples, $\text{tuples}(p)$.
- Every tuple $\langle v_1, \dots, v_k \rangle$ has k components and each components v_i is document vector over T
- Every tuple has a score, that is, $\text{score}(\langle v_1, \dots, v_k \rangle \in \text{tuples}(p))$
- In a base relation, this score is initialized to 1.



Conjunctive Queries

- Conjunctive query in WHIRL is written as $B_1 \wedge B_2 \wedge \dots \wedge B_k$ where each B_i is a literal.
- There are 2 types of literals

- EDB literal is written $p(X_1, \dots, X_k)$ where p is the name of relation and X is variable symbol.

The score of EDB literal will be :

$$\text{score}(B\theta) = \text{score}(\langle X_1\theta, \dots, X_k\theta \rangle \in p) \text{ if } \langle x_1\theta, \dots, x_k\theta \rangle \in p$$

$$\text{score}(B\theta) = 0 \quad \text{otherwise.}$$

- Similarity literal is written $X \sim Y$

The score of similarity literal will be :

$$\text{score}(B\theta) = \text{sim}(X\theta, Y\theta)$$

- We define the score for conjunction query as follow :

$$\text{score}(Q\theta) = \prod_{i=1}^n \text{score}(B_i\theta)$$



Union of Conjunctive Queries

- A basic WHIRL clause is written $p(X_1, \dots, X_k) \leftarrow Q$
- Union of queries in WHIRL is called view
- Support of view is defined as a set of triples $(A \leftarrow Q, \theta, s)$ satisfying conditions :
 - $(A \leftarrow Q) \in V$
 - $A\theta = a$ and $Q\theta$ is ground
 - $\text{score}(Q\theta) = s$ and $s > 0$
- Score of $\langle x_1, \dots, x_k \rangle$ is :

$$\text{score}(\langle X_1, \dots, X_K \rangle \in p) = 1 - \prod_{\langle C, \theta, s \rangle \in \text{support}(p(X_1, \dots, X_k))} (1 - s)$$



A* Search

```
procedure A*( $r, s_0, \text{goalState}(\cdot), \text{children}(\cdot)$ )
begin
  OPEN :=  $\{s_0\}$ 
  while (OPEN  $\neq \emptyset$ ) do
     $s := \text{argmax}_{s' \in \text{OPEN}} h(s')$ 
    OPEN := OPEN -  $\{s\}$ 
    if goalState( $s$ ) then
      output ( $s, h(s)$ )
      exit if  $r$  answers printed
    else
      OPEN := OPEN  $\cup$  children( $s$ )
    endif
  endwhile
end
```



Functions & States

- WHIRL uses 2 functions :
 - Inverted index : map terms $t \in T$ to the tuples that contain them
 - $\text{Index}(t, p, i)$: return tuples $\langle v_1, \dots, v_i, \dots, v_k \rangle$ in $\text{tuples}(p)$ in which $v_i^t > 0$
- States :
 - Pairs of $\langle \theta, E \rangle$ where θ = substitution and E = set of exclusion
 - E is a pair of $\langle t, Y \rangle$ where t = term, Y = variable
 - E is E -valid if it satisfies condition :
$$\forall \langle t, Y \rangle \in E, (Y\theta)^t = 0$$
 - Initial States: $\langle 0, 0 \rangle$
 - Goal States: $\langle \theta, E \rangle$ in which θ is ground.



Generating States

- Exploding States :
 - For a state $s = \langle \theta, E \rangle$, pick EDB literal $p(y_1, \dots, y_k)$ in which Y is unbound by θ .
 - Construct all states of the form $\langle \theta \cup \{y_1=v_1, \dots, y_k=v_k\}, E \rangle$ such that
 - $\langle v_1, \dots, v_k \rangle \in \text{tuples}(p)$
 - $\theta \cup \{y_1=v_1, \dots, y_k=v_k\}$ is E-valid
- Constraining a States :
 - For a state $s = \langle \theta, E \rangle$, pick constraining literal $X \sim Y$ and some term t (probably a term with highest weight) in document $X\theta$ such that $(t, Y) \in E$
 - Will produce 2 sets :
 - Singleton set containing state $s' = \langle \theta, E' \rangle$ where $E' = E \cup \{(t, Y)\}$
 - Set s_t contains all states $\langle \theta_i, E \rangle$ in which $\theta_i = \theta \cup \{y_1=v_1, \dots, y_k=v_k\}$ for $\langle v_1, \dots, v_k \rangle \in \text{index}(t, p, l)$ and θ_i is E-valid



Heuristic Function

- A* search algorithm uses heuristic function to determine which state has the highest score
- A heuristic function h is said to be admissible if for all states s and all states s' reachable from s , $h(s) \geq h(s')$
- Heuristic function is defined as follows :

$$h(\theta, E) = \prod_{i=1}^k h'(B_i, \theta, E)$$

Where :

$$h'(B_i, \theta, E) = \text{score}(B_i \theta) \quad \text{if } B_i \theta \text{ is ground, or}$$

$$h'((X \sim Y), \theta, E) = \sum_{t \in T: \langle t, Y \rangle \notin E} x^t \cdot \text{maxweight}(t, p, l) \text{ otherwise}$$



Steps to Generate States

1. Choose the smallest relation in the query
2. Do the explode state for the relation
3. Add each of the state to the OPEN list
4. Find the highest score state by using heuristic function
5. If the highest score state is ground, output it
6. Do the constraining state to the highest score state with constraining literal defined in the query
7. Repeat to step 4 until we have r-tuples in output



Example

- Query : $\text{const}(\text{IO}) \wedge \text{p}(\text{Company}, \text{Industry}) \wedge \text{Industry} \sim \text{IO}$
IO = “telecommunications services and/or equipment”
 - Start with the smallest relation, in this case IO
 - Explode IO, the result is single document IO, named it as state s_1 and placed it in OPEN list
 - Because s_1 is the only state in the OPEN list, removed it and do the constraining nodes to s_1
 - Choose term t which has the highest weight, probably “telecommunications”. The inverted index is used to find all tuples containing the words “telecommunications”

All of these substitutions are ground, so the value of $h(.)$ is the actual score of each tuple.
 - Put all of these tuples as states in OPEN list, together with new state s_1' containing the exclusion $\langle \text{telecommunications}, \text{Industry} \rangle$



Example (cont.)

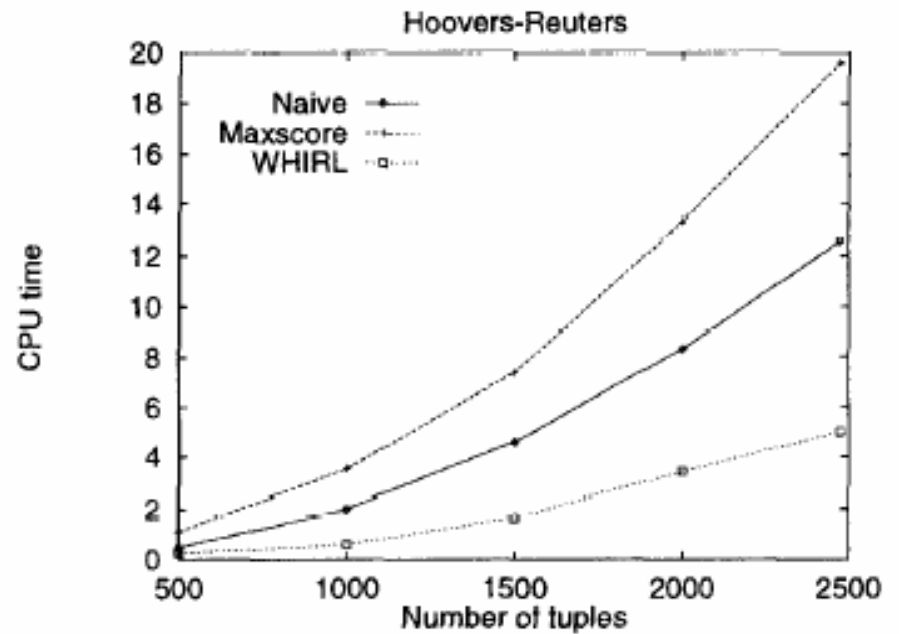
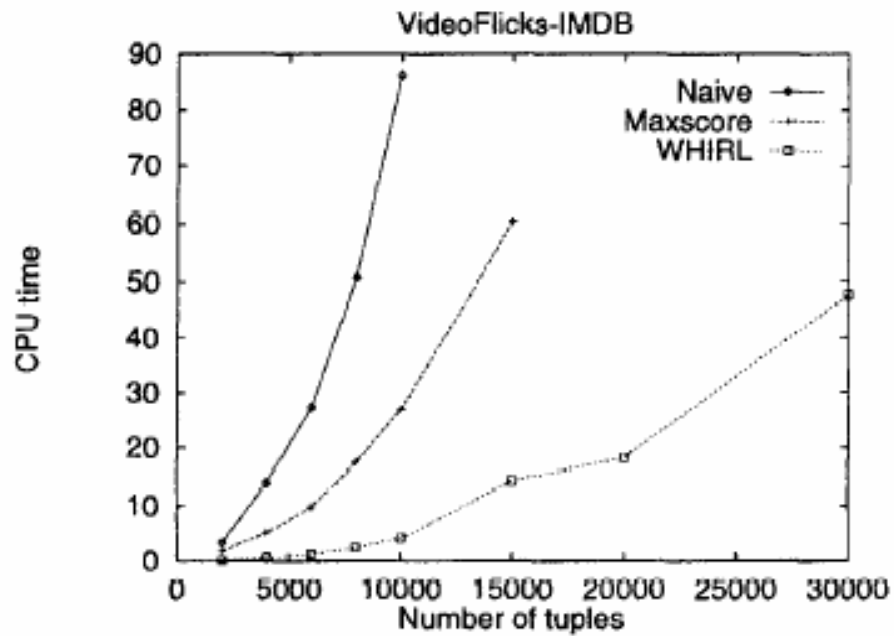
- Choose a state to be removed from the OPEN list
- Since we choose a term that has maximum weight, it might be that $h(s_1')$ is less than $h(.)$ from other states
- If it is, ground substitution will be removed from OPEN list, and an answer will be output
- If $h(s_1')$ is higher than the best goal statement we do again the constraining nodes to s_1' which means choosing a term from "services", "and", "or", "equipment" which has the highest weight
- Do like what have we done before, until we found r-answer as the output



Experimentation

- The experimentation compares WHIRL with 2 other methods
- Naïve method for similarity join
 - Take i -th column of relation p , submit it as IR ranked retrieval query to a corpus corresponding to the j -column of relation q
 - Top r result is merged to find the best r -answer
- Maxscore optimization
 - The same for naïve method, except that maxscore optimization is used to find best r result.

Time Result





Average Precision I

- The average precision is computed as follow :

$$averageprecision = \sum_{k=1}^r c(k) \cdot \frac{a_k}{k}$$

Where : a_k = number of correct answer in first k

$c(k)$ = 1 if the k -th answer is correct and 0 otherwise

Similarity Joins

Domain	Average Precision
Business	84.6%
Animals	92.1%
Movies	100.0%



Average Precision II

Similarity Joins with
Incompatible Schemata

Pairing	Average Precision
movieName/movieName	100.0%
movieListing/movieName	100.0%
movieName/review	98.0%
movieListing/review	94.0%



Conclusion
