

**Task 1: Query Hierarchy**

Given is the following database, with two additional tables: “Cite” (to join Paper and Paper) and “Write” (to join Author and Paper).

Author	
id	name
1	Charlie Carpenter
2	Michael Richardson
3	Michelle

Paper	
id	title
1	Contributions of Michelle
2	Keyword Search in XML
3	Pattern Matching in XML
4	Algorithms for TopK Query

1.1 Create the query hierarchy for the query  $K = \{Michelle, XML\}$  where no query has more than two joins.

1.2 Which query construction option would you definitely not present to the user?

**Task 2: Estimating the Probability of a Query Interpretation**

Given are the following database, keyword query and query:

Actor	
id	name
1	Tom Hanks
2	Collin Hanks
3	Tom Cruise

Movie		
id	title	year
1	Catch Me If You Can	2002
2	Cast Away - Verschollen	2000
3	Scene by Scene	2001

$$K = \{hanks, 2001, cruise\},$$

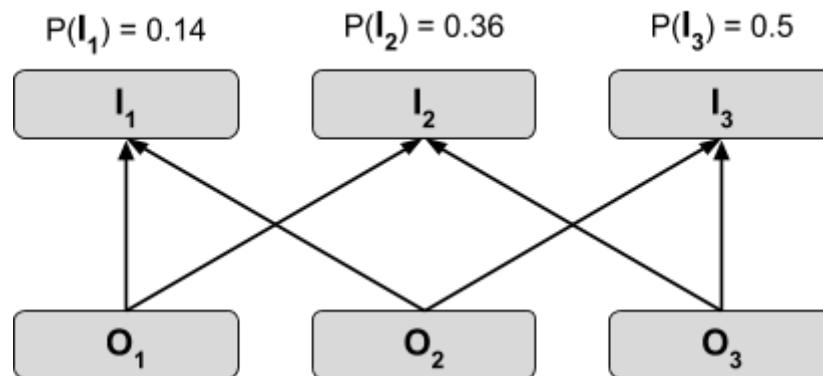
$$Q = \sigma_{hanks \in name}(Actor) \bowtie Acts \bowtie \sigma_{2001 \in year}(Movie) \bowtie Acts \bowtie \sigma_{cruise \in name}(Actor).$$

Task 2.1 Explain in words what  $P(Q|K)$  describes.

Task 2.2 Estimate  $P(Q|K)$ .

### **Task 3: Selecting a Query Construction Option**

Given is the following query hierarchy:



Predict which query construction option ( $O_1, O_2$  or  $O_3$ ) should be presented to the user.

Task 2.1 Explain in words what the graph shows and which criteria lead to the QCO selection.

Task 2.2 What do you need to compute?

Task 2.3 Compute the needed values and decide for a QCO.

## Appendix: Formulas

QCP: Query Construction Plan  
QCO: Query Construction Option  
IG: Information Gain

### Probability of a query interpretation

$$P(Q|K) = P(I, T|T) \propto \left( \prod_{k_i \in K} P(A_i : k_i|A_i) \right) \cdot P(T)$$

$I$  – the set of keyword interpretations  $\{A_i : k_i\}$  in  $Q$

$T$  – the template of  $Q$

### Probability of a keyword interpretation

$P(\sigma_{k_i} \in A_i | \sigma_q \in A_i)$  can be estimated using Attribute Term Frequency (ATF):

$$ATF(k_i, A_i) = \frac{TF(k_i, A_i) + \alpha}{N_{A_i} + \alpha \cdot B} \quad \text{– the normalized keyword frequency of } k_i \text{ in } A_i$$

$N_{A_i}$  – the number of keywords in  $A_i$

$\alpha$  – a smoothing parameter, typically  $\alpha = 1$  (Laplace smoothing)

$B$  – the vocabulary size

### Probability of a query template

$$P(T) = \frac{\#occurrences(T) + \alpha}{N + \alpha \cdot B}$$

$\#occurrences(T)$  – number of queries in the log using  $T$  as a template

$N$  – the total number of queries in the log

$\alpha$  – a smoothing parameter, typically set to 1 (Laplace smoothing)

$B$  – the vocabulary size

When the query log is absent or is not sufficient, we assume that all query templates are equally probable.

### A measure of QCO efficiency and probability estimation for QCOs

$H(\zeta) = -\sum_{I \in \zeta} P(I) \cdot \log_2 P(I)$  – entropy of the query interpretation space

$$IG(O) = H(\zeta) - H(\zeta|O)$$

– the expected information gain of a QCO as entropy reduction

$$H(\zeta|O) = P(O) \cdot H(\zeta_{|O}) + P(\neg O) \cdot H(\zeta_{\neg O})$$

– the entropy of the interpretation space given the QCO

$$P(O) = \sum_{I \in \zeta(O)} P(I)$$

– the probability of a QCO using probabilities of the subsumed query interpretations