



(1) The incorporation of language-model-based frameworks to resource descriptor allows us to calculate the similarity between collection or collection and queries based on more solid theoretical ground.

(2) A topology reorganization protocol allowing semantically similar nodes together to form loose content clusters in distributed manner. We then analyze the impact of topological factors on IR efficiency.

(3) Two efficient query routing algorithms which take advantage of the language model and the underlying topology.

We evaluated these algorithms on REC 100 and REC VLC (Very Large Collection 1) respectively [7]. The large collection is split to hundred of small sub-collections by



At this point,  $n_i$ 's routing table is reorganized according to the following rule:

( )  $\alpha$ % of its degree is redesigned to its most similar neighbor while the rest  $(1-\alpha)$ % neighbors are randomly picked from the set  $\{n_j, r_i, \bigcup_{n' \in n_j, r_i} \{n_p, r_n\}$ . Choosing neighbors randomly in an effort to keep the graph well connected. An experiment shows that if all the neighbors are chosen from the most similar node, the resultant network suffers from bad connectivity. Specifically, it contains many

(KL) divergence to measure the distance between collection model or collection model and query model. The formula is:

$$D(p||q) = \sum_i p(v) \log \frac{p(v)}{q(v)}$$

Unfortunately, this formula includes every possible word, so it is very time-consuming. To speed the process, an approximate formula is used: [9]

$$D(p||q) = - \sum_{w \in V \cap Q} p(w) \log \theta_Q$$