

# Improving Unsupervised Query Segmentation using Parts-of-Speech Sequence Information

Rishiraj Saha Roy, Yogarshi Vyas, Niloy Ganguly (IIT Kharagpur)  
Monojit Choudhury (Microsoft Research)

**Query segmentation** is the process of breaking down Web search queries into their constituent structural units, thereby aiding the retrieval process. This study connects two orthogonal approaches to segmentation or chunking of text fragments – those that rely on purely statistical word association measures and those that try to incorporate linguistic information, used commonly for Natural Language chunking.

Our initial experiments show that POS tagging does improve query segmentation. Although we do not observe any performance benefits of using a specific POS tagset or tagging approach over others, we do observe that tagset and taggers designed for POS tagging English text help improve the segmentation of a complementary set of queries than the ones which are benefitted by unsupervised POS induction. Thus, appropriately combining these two approaches is expected to lead to better segmentation.

## Example Clusters from Bie-S

Cluster 1: bake, casserole, dessert, fry, meatballs,  
Cluster 2: athletics, baseball, cycling, football,  
Cluster 3: army, citizenship, customs, defence,  
Cluster 4: battlefield, diablo, godfather, hitman,

