**Clustering and Summarizing Medline Abstracts**
*Dunlavy, Daniel M.*[*1], *Conroy, John M.*[2], *O'Leary, Timothy J.*[3], *O'Leary, Dianne P.*[1]
[1]*University of Maryland, College Park, MD, USA;* [2]*Center for Computing Sciences, Bowie, MD, USA;*
[3]*Armed Forces Institute of Pathology, Rockville, MD, USA*

The amount of reference material available to medical professionals is overwhelming, and its usefulness depends not only on timely delivery but also on information triage. The MEDLINE system is invaluable for extracting documents relevant to a query and presenting them in a rank-ordered list. When a query returns a large number of documents, however, it is often impractical to search this list for desired documents.

In previous work, we have developed a system to extract documents relevant to a query, cluster these documents by subject, and return a summary of the documents in each cluster augmented by the list of documents in that cluster. This Query-Cluster-Summarize (QCS) system was designed with newswire documents in mind. We have tested it on MEDLINE documents related to gastrointestinal stromal tumors (GIST), and the results are encouraging.

A screenshot of the QCS system is shown below for the query "kit immunohistochemistry". The interface to the QCS system includes 1) an *input* section for the query and choice of document set, 2) a *navigation* section with links to clustered documents (Q: top documents retrieved for the query and their scores, C: documents from which the summary sentences were drawn and the sentence indices), and 3) an *output viewing* section, which here contains the default output of multiple document summaries for the first two topic clusters.