

# Derivation of Knowledge from Digital Content

A. Shelly Spearing, Jorge H. Román, HPC-1; M. Linn Collins, STBPO-RL

As electronic content proliferates, it becomes nearly impossible to fully consume and assess all of the available information. Over the past 7 years, LANL's Digital Knowledge Factory (DKF) has created a suite of digital-content-analysis tools. This combination of commercial off the shelf (COTS), Open Source, and homegrown code is used to gather, reduce, annotate, organize, synthesize, and visualize digital content for human consumption.

The tools can be applied to collections of text-based documents from virtually any source. The algorithms go beyond traditional natural language processing and statistical analysis—word-location algorithms automatically extract the gist of the content, while others annotate targeted concepts, organize documents, and calculate goodness-of-fit with respect to a specified conceptual area. Additional modules extract features, such as dates and locations, and group documents for comparative analyses. In some cases we are also able to compute the trustworthiness and/or mood of the author. When looking at larger collections, we categorize subject-matter expertise, emerging and fading trends, and distill entire collections into a variety of single-page graphical representations.

Structured information (i.e., metadata) can augment the digital knowledge to facilitate analyses of time trends, geographical colocation, and authorship, among others. Through information reduction, annotation, fusion, and organization, the analyst is able to assimilate content and form hypotheses more quickly.

A goal of the DKF project is to expedite knowledge assimilation by synthesizing digital content into a

set of knowledge visualization schemes. In particular, it is hoped that an information consumer can look first at a few graphical representations of the concepts contained within thousands of pages of text, draw conclusions about the documents in aggregate, formulate hypotheses, and then focus attention on the particular documents that are relevant to the conclusions and hypotheses. Through a reduction process we focus the analyst on the important concepts and the relationships among them.

A first step in trend identification is depicted in Fig. 1, a timeline summary of 475 documents containing “Iraq” that were extracted from the White House Press Archive website. DKF tools were used to identify the top-level concepts

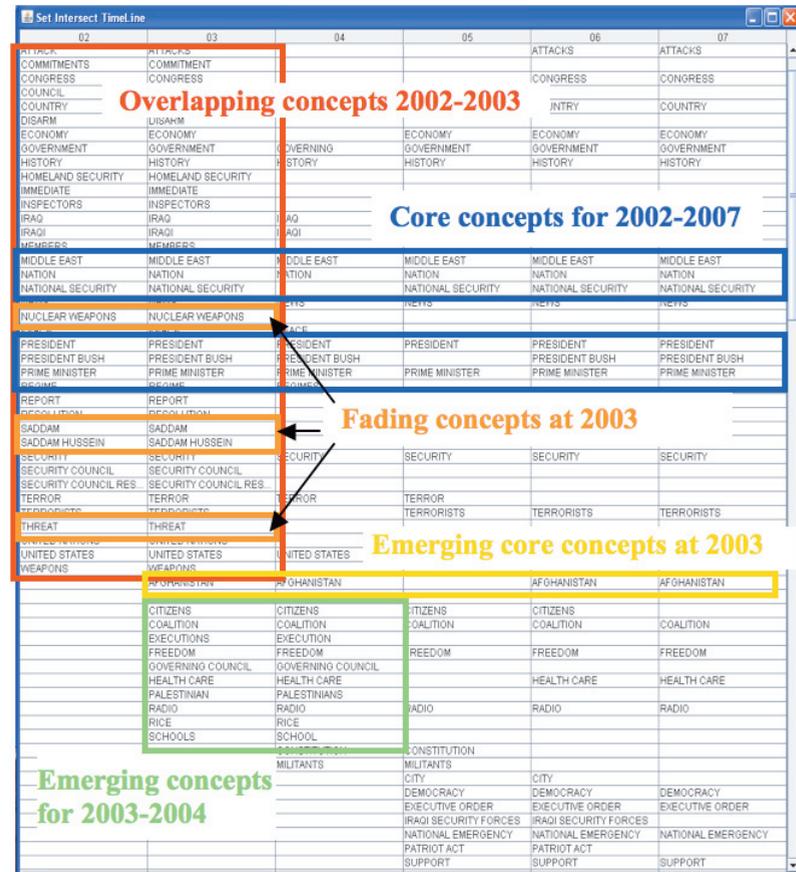


Fig. 1. Timeline summary (right) of 475 documents.

