

The New York Times Here we found numerous duplicated branches in the constructed DAG (e.g., for research, television, education, medicine, defence and military forces). This indicates that for these topics, two distinct sets of tags were used in parallel. The DAG is much less organised than that of the Spiegel and of the Guardian. There are 31 isolated components, most of them correspond to one theme (e.g. "Baseball"). The sizes of the components varies from 898 to 2, and there is a continuous range of them from the 2nd largest one (274 tags) down. There are no very general categories. Although a number of large related components exists (under the tags "Basketball", "Baseball", "Football"), these components are not collected under a general "Sport" tag. It seems as if there were no demand for using general tags. Note that there is a tag called "sports", however, it appears only on 5 news items, and it is negligible. A technical consequence is that the DAG construction algorithm does not always select the most general tags as roots, because they lack the important connections to other components. Instead, one of the more specific tags can be selected for a central position, for example, "Middle East and North Africa Unrest (2010-)" for foreign affairs, or "European Sovereign Debt Crisis (2010-)" for Europe-related tags. In other words, the centrality no longer correlates only with the generality for the top tags. Some lower-level branches end up at unexpected places, e.g., **Environment** under **Iran**. Superfluous levels appears, for example, **International Relations** under **United States International Relations**.

The Australian The DAG looks disorganised overall. There are about 1900 components for the 79504 tags without the pre-filtering, and about 300 components for the min. 5 news items-filtered 1673. There are no macroscopic components, the largest one's size is just 3480 (out of 79504 tags) and 165 (out of 1673 tags), which is less than 10% of the total nodes. Even the existing components look more like just bunches of more or less associated tags than small hierarchical structures.

In general, the top of the constructed DAGs are much better than the bottom. This is no surprise - there is much more information for the construction algorithm at the top of the DAG.

Pairwise comparisons

We carried out a pairwise comparison between the journals from the point of view of their content organisation. Since the audience and the interests of the journals are different, the list of tags appearing on the articles was unique for each news portal. Therefore, before actually comparing the tag hierarchies, first we needed to create a common tag set for each pair of journals. In a number of cases, finding the corresponding tag pairs went beyond a simple string matching and was based on semantic matching, e.g., "Fossil fuels" (Guardian) was matched with "Oil