# DiversiNews: Surfacing Diversity in Online News

*Mitja Trampuš Flavio Fuart, Daniele Pighin, Tadej Štajner, Jan Berčič, Blaž Novak, Delia Rusu, Luka Stopar, Marko Grobelnik*

■ *For most events of at least moderate significance, there are likely tens, often hundreds or thousands of online articles reporting on it, each from a slightly different perspective. If we want to understand an event in depth, from multiple perspectives, we need to aggregate multiple sources and understand the relations between them. However, current news aggregators do not offer this kind of functionality. As a step toward a solution, we propose DiversiNews, a real-time news aggregation and exploration platfom whose main feature is a novel set of controls that allow users to contrast reports of a selected event based on topical emphases, sentiment differences, and/or publisher geolocation. News events are presented in the form of a ranked list of articles pertaining to the event and an automatically generated summary. Both the ranking and the summary are interactive and respond in real time to user's change of controls. We validated the concept and the user interface through user tests with positive results.*

The Internet has been strongly gaining prominence as a news medium; in 2012, it overtook the TV as the most used news source for people under 30 in the United States (Kohut 2013). In addition, it has significantly changed the way in which many people find and consume news. Multiple publishers are now reachable more easily than ever before. Social bookmarking sites present us with news deemed interesting by our peers. News aggregation sites give us an instant overview of the topics of the day.

Although this plethora of sources theoretically provides a richness of information that even 15 years ago was unthinkable, practice can prove it much harder to find multiple and truly varied views on a given subject. Consider the following fictional but reasonable examples.

Alice is browsing the Internet when she encounters an article saying that Coca-Cola announced a new shape for its bottle, a first in many years. Since Alice owns some Coca-Cola stock, she

is curious to know more, especially about the likely business implications. It turns out that the general public is primarily interested in the history of Coca-Cola bottle design and after searching online, Alice finds mostly articles on that topic as those are the most popular. The comments on Reddit are similarly narrow in focus: the ones about the design are popular, upvoted, and displayed prominently.

Bob is interested in politics and would like to know more about the developing civil unrest in Elbonia. He is savvier than Alice and uses Google News to efficiently obtain a large number of reports on the issue. In fact, there are literally hundreds of reports and Bob is overwhelmed — he has no easy way of finding and contrasting the leftist and the rightist opinions, the local or international points of view, the articles in support of the rebels, and the progovernment ones.

There is nothing special about Reddit or Google News in this context; the above anecdotes could just as easily take place on any major news sharing or aggregation website. They serve to illustrate broader, general issues that are being faced daily by users browsing news on the Internet: a single point of view, and information overload.

*Single Point of View:* A single article almost always means a single author, limited perspectives, and only partial coverage of an event. To access multiple points of view, we may turn to news aggregation sites like like Google News, Bing News, NewsVine, or others. However, even these sites represent each event with only one or maybe a handful of articles. There is a clear incentive to promote the most popular articles, thus making them even more popular and consequently exposed; a classic rich-get-richer scheme, evidenced by for example, Munson, Zhou, and Resnick (2009); Brzozowski, Hogg, and Szabo (2008); and Leskovec, Kleinberg, and Faloutsos (2005). These sites optimize for discoverability of events, not diversity of coverage.

A similar effect happens on social link sharing websites (Facebook, Reddit, Pinterest, Flipboard, and others). The promote — upvote self-fulfilling cycle gives rise to a single way of thinking, pushing fringe opinions and content further into obscurity. Current news sites do not provide an easy way to surface the diversity in the data.

When reader come across an article on a new event, they often do not have the necessary contextual knowledge that would allow them to put that piece into perspective and judge the novel information. Such questions can be answered by providing access to the background on the topic, the involved people and organizations, the places affected by the event, and where that article fits into the overall opinion spectrum.

*Information Overload:* While existing news aggregators are reasonably good at collecting large amounts of articles reporting on a single news story or issue, users are mostly left to their own devices when it comes to navigating those articles. Typically, we can filter or sort by relevance and time. However, articles on a single issue can differ in much more: their provenance, trustworthiness, fact coverage, topical focus, point of view, and so on. Current news sites provide no way of navigating according to these criteria.

In short, the diversity in news reporting is underrepresented on the Internet. In addition, individual news sources are reducing the amount of editorials and commentary, while simultaneously, people of each coming generation spend less time reading the news even as they age, according to a Pew[1] study (Kohut 2013). In other words, diverse views are becoming scarcer and people are willing to invest less and less time into finding them. This is clearly an undesirable situation that we should fight against: it has long been known that informed citizens lead to better policy choices and thus better life overall (Downs 1957, Milner 2002).

In this article, we propose a software system that collects news articles by crawling the Internet, groups them into stories (that is, clusters of articles reporting on the same event or issue), and presents them through a novel user interface that helps readers discover contrasting perspectives on the news. The central screen of the application, focusing on a single story, presents an overview of the contributing articles: what aspects of the story they emphasize, where in the world they were written and whether they view the story in a positive or negative light. The individual articles are also presented, along with an automatic summary of the story. The user can reorder the articles based on any combination of the modalities previously mentioned (subtopic or aspect, geography, sentiment) to surface a specific point of view. The summary also changes to reflect the new focus of interest.

We developed two front ends: a prototype version for the web (figure 1) called DiversiNews, and a more robust and user-friendly one for iOS (figure 2), called iDiversiNews. An extensive evaluation, carried out on the interface of the web service, provided useful feedback that we incorporated in the more refined iOS application.

We claim two main contributions: functionality and user interface.

*Functionality:* We collect publicly available articles in real time and infer attributes like (sub)topical focus, sentiment, and place of origin. Through this, we enable users to explore news in a manner significantly different from existing tools, with much easier and more explicit access to the diverse views on a topic.

*User Interface:* We synthesize information across multiple modalities (geography, topics, sentiment) by consistently mapping each of them onto a two-dimensional surface area on the screen, with continuous movement of controls corresponding to continuous changes in each of the modalities. This

*Figure 1. The Prototype Web Interface for Surfacing Diversity in News.*

representation is particularly well suited to touch-screen interfaces. Also novel is the near real-time adaptive summary: users get immediate feedback on their selection of perspective, without having to read several articles.

The full-fledged iPad version is available in the iTunes app store.[2] The prototype web version is available online.[3]

## Scope

Not everyone has a need for diversified browsing, and certainly not all the time. For some topics, news content does not vary much across sources simply by nature, for example, because it is highly technical. On top of that, people in general prefer to keep

returning to a handful of information sources, regardless of what else is available (Flaxman, Goel, and Rao 2013). The ecosystem of online news is not getting more polarized with time (Gentzkow and Shapiro 2011) — we seek out our bubbles a lot of the time.

However, this is far from saying there are no situations in which discoverability of diverse and rich content is desirable. Alice and Bob from the introductory examples, a fan of Michael Jackson reading about the performer's death, a public relations (PR) person following media reports on his or her company's big public event, an active citizen trying to keep up to date with a proposed policy change, and so on. With DiversiNews, we aim to serve these and similar

needs, where a reader wants or needs to understand a topic in-depth, but existing tools fail easily to provide multiple views or help understand the relations between them. We do not limit ourselves to a specific subset of users or topics, although some of them are more likely to be the subject matter or audience of DiversiNews than others. For example, we expect political topics are more amenable to subjective writing and differing opinions, and technical topics less so; and that specialists and students are more likely to need to understand a story in detail.

Going beyond anecdotal use cases, a recent user study by Munson and Resnick (2010) evaluated what happens if the diversity is forced upon (or away from) the user. Test subjects were asked about their political preferences and then exposed to a collection of news that agreed with their preferences to varying extents. Two groups of users were discernible: one was happiest if all the articles agreed with their views, while the other was happiest when served a balanced mixture of news that both support and challenge their views. Although these users represented a minority, they indicate there is a target audience for technologies that make diverse content more accessible.

# Related Work

In related work, the phenomenon of readers' exposure to a limited set of viewpoints is called media bias, selective exposure, or filter bubble — the names differ in nuances but are closely related. The opposite, that is, access to a range of viewpoints, is quite uniformly called diversity, a term used in this article as well.

## The Case for Diversified Browsing

There is a large body of research associated with identifying, measuring, and explaining media bias. Frequently, the research in this area focuses on diversity and biases along a single dimension, typically the political orientation (liberal versus conservative). For example, An, Quercia, and Crowcroft (2013a) tracked Facebook users' patterns of sharing links to articles and confirmed that liberals were much more likely to share liberally inclined articles and vice versa for conservatives.

Maier (2005) surveyed several thousand news sources cited in newspapers and found factual or subjective disagreement between the sources and the citing articles in 61 percent of the articles. This shows that in order to get objective information, one should ideally have easy access to multiple articles on a story.

Voakes and Kapfer (1996) analyzed multiple news stories and found that the content diversity is on average substantially lower than the source diversity; in other words, simply reading a high number of sources does not necessarily provide diverse content. This suggests that diversity-aware news browsing systems should understand news on some level, be aware of its content and other attributes.

## News Recommender Systems

News recommender systems stand to benefit from similar data representations and similarity metrics as our work. Among other things, they aim for serendipity — exposing users to content that is relevant but different enough from what they've already seen, so that they can discover new stories.

An extensive survey of news recommenders was published by Borges and Lorena (2010). Standard news recommender systems model the user and recommend new articles based on content similarity (content-based filtering) or based on past actions of similar users (collaborative filtering).

Like ourselves, Tavakolifard and colleagues (2013) use BOW (bag of words) and geography as features, and additionally suggest using other named entities. Lv and colleagues (2011) use standard BOW features but suggest new similarity metrics that better capture, for example, connection clarity.

## Diversity-Aware News Browsing

While the work listed is mostly descriptive in nature, there is also no lack of prescriptive research trying to provide solutions that would ameliorate the current state of affairs. In his Ph.D. thesis, Munson (2012) suggests several visualizations of a user's browsing patterns, for example a graph of the prevalence of liberal-leaning articles among those read by the user. As the graph evolves through time, the user can track her reading habits, holding herself accountable to a balanced diet of opinions. This complements our work where the goal is not to identify a user's need for balanced reporting, but rather to help her satisfy that need.

Very closely related to our work is NewsCube by Park et al. (Park, Lee, and Song 2010; Park et al. 2012), a system for news aggregation, processing and diversity-aware delivery. DiversiNews and NewsCube have much common — they both choose to offer diversity through a stand-alone news portal, and a lot of the preprocessing work is therefore similar across the two systems. There are however notable differences in delivery. For one, NewsCube offers no interactive exploration but rather groups and ranks articles within a story in a fixed way that is hoped to offer maximally diverse information in one screenful. Secondly, NewsCube focuses on topical diversity only.

Later work by the same authors extends the information presented by NewsCube with a more detailed characterization of biases and a novel data acquisition method. NewsCube 2.0 (Park et al. 2011) is a browser add-on that allows users to collaboratively tag articles with the types of exhibited biases (for example, omission of information, suggestive photo, subjective phrasing, and others) and place them on the *framing spectrum,* that is, decide how strongly lib-

eral or conservative the article's outlook is. User input is then presented in the NewsCube interface.

Also similar in spirit is BLEWS (Gamon et al. 2008), a news browser that annotates news articles with color bars representing the amount and affect (positive versus negative) of links to that article from liberal and conservative blogs. The political affiliation of blogs is taken from a third-party list; the affect is machine learned. As an almost-inverse of BLEWS, the hobby project Memeorandum Colors[4] by Andy Baio provides a Firefox plug-in that colors links to bloggers according to their political leaning (red versus blue). The political affiliation of blogs is calculated with SVD, each blog being originally represented with a list of news articles to which it links.

Another noteworthy and much more mature news portal is the Europe Media Monitor (Steinberger, Pouliquen, and Goot 2009), which aims to bring together viewpoints across languages. The website offers a number of news aggregation and analysis tools that track stories across time, languages, and geographic locations. It also detects breaking news stories and hottest news topics. Topic-specific processing is used, for example, to monitor EU policy areas[5] and possible disease outbreaks (Rortais et al. 2010).

In a similar vein, DisputeFinder (Ennals, Trushkowsky, and Agosta 2010) is a browser extension that lets users mark up disputable claims on web pages and point to claims to the contrary. The benefit comes from the collaborative nature of the tool: when browsing, the extension highlights known disputed claims and presents to the user a list of articles that support a different point of view.

In contrast to most of the work that focuses on political diversity, Zhang, Kawai, and Kumamoto (2012) identified similar and diverse news sources in terms of the prevalent emotions they convey.

## Diversity in Other News Modalities

News in the traditional form of articles is among the most amenable to analysis. For news in other forms (video, tweets), the promotion of diversity is mostly restricted to attempts at making the data collections more easily navigable.

Social Mention[6] is a social media search and analysis platform that aggregates different user generated content, providing it as a single information stream. The platform provides sentiment (positive, negative, and neutral), top keywords, top users, or hashtags related to the aggregated content.

The Global Twitter Heartbeat (Leetaru, Wang, and Cao 2013) project performs real-time Twitter stream processing, based on a 10 percent sample of the Twitter feed. The text of each tweet is analyzed in order to assign its location. A heat map infographic displays the tweet location, intensity, and tone.

## Text-Mining Methods

Finally, many of the data processing components that iDiversiNews and DiversiNews uses behind the scenes and combines in a novel way have been researched before in an isolated setting. It is however beyond the scope of this article to provide an extended technical overview of existing work on each of the components.

# User Interface

Traditionally, publishers and news aggregation services create a particular, static view on a news story. Our guiding principle was that there is no default story overview and that no aggregation fits all purposes.

An important consideration when designing the user interface was to allow users to navigate and explore different modalities of a story. The challenge here is to show the big picture, thus reducing information overload, but still allowing the drill-down to raw news articles. The latter is very important to strengthen the trustworthiness of the system: at every step, users should be allowed to verify the original content that contributed to the aggregated view created by the system.

The iOS application consists of five screens, depicted in figure 2. The usage flow is linear: you start at the screen in figure 2a, select a story at figure 2b, and analyze and manipulate it at figure 2c, the main screen, which shows a single story at a glance. In case you require more information on the story, you can continue to screens figure 2d and figure 2e. We describe each of the screens in the subsections.

The web version, which served as a prototype for the iPad application, is somewhat simpler and consists of a story exploration screen (figure 1) that is functionally equivalent to the main iOS screen (figure 2c), and a very limited story discovery screen listing a fixed set of test stories.

## Story Discovery

When you start using the application, you are presented with the screen in figure 2a. Your first task is to find a story into which you want to get a deep insight from multiple viewpoints. There are several ways you can find such a story. You can choose to list stories pertaining to the most active DMOZ[7] categories (see the Topic Classification and Entity Extraction subsection for technical details), the most active publishing locations (see the Data Aggregation subsection), or simply all the recent stories. A list of stories that match your selection appears on the next screen (figure 2b). You can further filter news stories by searching for keywords (see the Data Aggregation subsection) and selecting world regions (see the Geotagging subsection) or topics.

*Figure 2. The Screens in iDiversiNews.*

(a), (b) Story discovery, where you find and select a story of interest. (c) Story exploration for a single story, the most important screen of the application, with controls for surfacing diversity. (d) Article list for a single story, with articles sorted according to controls in (c); and e) Detailed article view of a single article and its metadata.

*Figure 3. System Architecture.*

## Story Exploration

The most prominent and novel part of the application is the story exploration screen (large screenshot in figure 4) where different views on a single story can be explored. The screen is conceptually divided into thirds. The top third provides metadata about the story. The middle third offers controls for expressing interest in articles in terms of publisher location, subtopic, or sentiment. Based on the expressed interest, the bottom third is regenerated on the fly: it contains an automatically generated summary (on the left) and the most relevant articles (right).

The user controls work as follows: The Publisher Location control (figure 4 marker C) shows the indi-

vidual articles on a world map, along with a target icon. By moving the target to any location on the map, you instruct the system to focus on articles published in that part of the world.

The Subtopic control (figure 4 marker D) is the least conventional in appearance. It shows a map of subtopics, represented by keywords. Similar subtopics and keywords are displayed more closely together. Again, you are presented with a target icon that you can move to any part of the topic space to focus on that subtopic. See the Subtopic Detection subsection for technical details.

The Sentiment control (figure 4 marker E) is a simple slider that spans the sentiment spectrum of the articles and allows you to focus on articles with a positive or negative tone. For example, moving the slider to the far left gives you the most negative articles (or the least positive ones in the case of inherently positive news).

Both maps, the geographical and topical, serve the double function of giving you an overview of the active areas (in either the geographical or topical space) but also allowing you to focus on an area.

### Design Considerations

In choosing the dimensions users can navigate, we worked under the restriction of fitting everything on one screen, a necessity for making the discovery process fast and truly interactive. Equally important, the complexity of interaction had to remain manageable. With these considerations in mind, we settled on the three attributes described above — they are relevant to the task at hand, readily understandable, and can be reasonably fit onto a small screen.

The single-screen requirement however still poses interesting challenges to presenting all the information; these are exacerbated by the stark mismatch of the two-dimensional nature of the screen and the high-dimensional nature of text data. The sentiment is the least problematic in this regard; its one-dimensional value can be directly presented on the screen. Publishers' locations like Paris or the United States can be relatively easily merged with background knowledge bases to produce coordinates and intuitively displayed on the map; to accommodate the small space, heat-map-like summarization is required. Lastly, the inherently high-dimensional space of topics is the most problematic and required the most tailored solution, a combination of dimensionality reduction (two-dimensional projection) and data filtering (reduction to key words) to produce the topic map. Similarly cumbersome is the representation of the output, that is, the relevant articles. Here, we opted for data filtering alone, in the form of an adaptive summary and top article titles.

### Detailed Article View

If the Story Overview screen does not provide sufficiently detailed information, you can drill down your analysis to the level of individual articles by clicking any article in the lower right corner of that screen, or in the full Article List (figure 2d). This opens the article in the Detailed view (figure 2e), where you can browse the original article web page. Above it, we show metadata extracted by the system: publishing date, keywords, key people and organizations, sentiment, and so on.

## Data Pipeline

In this section, we describe the implementation and deployment of the system in overview. The forward references in the previous section point to most of the subsections that follow here, giving a clue as to how the individual components contribute to the system's functionality.

### Overall System Architecture

The chosen architecture (figure 3) is designed for scalability and reliability in an environment that allows near real-time processing, browsing, and analysis of news. It is service oriented to make it stable, easily maintainable, and extendable.

To enhance user experience, the iDiversiNews server[8] aims to preprocess as much data as possible before serving it to the client application: data is transferred from the Newsfeed (top of figure 3), a real-time feed of news articles, and processed as soon as it is made available. News articles are first enriched using the Enrycher pipeline (Štajner et al. 2010) (right-hand part of figure 3), a series of sequentially executed text annotation services: part-of-speech (POS) tagging, sentiment detection, named entity recognition (NER) and resolution, and categorization. In this part of the pipeline, articles are processed individually; in the later stages, they are grouped into stories and analyzed collectively.

The fully preprocessed articles are passed to the clustering service, thus creating news stories. Each story is then annotated with the following: region, country, keywords, categories, publishers locations, sentiment information, key entities, information on projection of the topical space into two-dimensions (see the Subtopic Detection subsection), and a representative thumbnail image. Most of this information is derived from the contributing articles using simple majority voting. In addition, the ranking of articles (see the Article Ranking subsection) and the corresponding summary (Summarization subsection) are calculated for the default position of user controls. All the information on a cluster is then stored into a Lucene-indexed repository.

At the end of the pipeline, a list of the most prominent stories in the news is periodically published through the Web API module and an aggregation of basic statistics is created. The client application can request the latest list of current news or trigger a customized search.

*Figure 4. A Closer Look at the Story Discovery and Story Exploration Screens.*

*Left:* Story discovery. *Right:* Story explanation. The hand symbols and the associated gesture indicators show possible user interactions and are not part of the interface.

## Data Aggregation

We use the Newsfeed (Trampuš and Novak 2012) system to collect the articles and their metadata. Newsfeed crawls RSS feeds of publishers across the globe and provides a real-time unified feed of news articles (figure 5) at a rate of about 100,000 per day. We convert all articles to plain text and detect their language. For details, please see the reference paper. In addition, Newsfeed provides information on publishers' geolocation (see the Geotagging subsection). About 50 percent of the content is in English; we currently discard non-English articles in DiversiNews.

## Article Clustering

To identify the news stories, we have to group articles according to events they are describing. This is achieved by using clustering algorithms. Clustering is a well established area of research, but traditionally it has not been developed for clustering stream-based data with a temporal dimension. Lately a number of approaches have been proposed (Silva et al. 2013).

## A Step-by-Step Example

Figure 4 illustrates a sample usage session by Shawn,

a (fictional) political activist. On the first screen, Shawn sees that *government* is a trending topic. His interest piqued, he selects it (action A in the figure). The next screen presents him with government-related stories; he selects the one about John Kerry trying for peace in the Middle East (action B). This brings him to the rightmost screen in the figure where he can explore the different view on the story.

First, Shawn wants to see how the Americans see the story, versus the Middle Eastern media support for Kerry's actions. By focusing the Publisher Location control (action C) first on the United States and then the Middle East, he discovers that the first view Kerry's actions as commendable if somewhat foolish and stress Kerry's opinion that the peace talks have to move forward; the Israeli media on the other hand is much more antagonistic, feels the peace solution is being forced on them, and emphasizes Kerry's veiled threats of a possible political boycott in case of failed negotiations.

Next, Shawn wants to understand the different aspects of the story. Using the Subtopic control (action D), he identifies the main aspects and selects them to find articles that focus on questions like "What is Abbas planning?" and "What exactly constitutes the boycott Kerry threatened to Israel?" and can give more detailed answers.

Finally, with the Sentiment slider (action E), Shawn can find articles that highlight the negative aspects, again mostly about the deep offense that Israel has taken with boycott threats. He can also see the positive articles that still cover the problems but focus more on the economic and humanitarian implications of a potential successful resolution.

In real time, as Shawn is making his selections, articles in the lower right corner are reordered accordingly and a corresponding summary of the most relevant articles is displayed in the lower left corner. All of these analyses would be much harder if Shawn were only presented with a list of all related articles ordered by time (which some news aggregators offer), and even more so if he tried to search for individual related articles online.

We have developed our system based on the algorithm by Azzopardi and Staff (2012). The key idea is to add each new article to the story most similar to it, or seed a new story with that article if all similarities to existing stories are below a threshold. Some additional details had to be solved and parameters tuned to our dataset. In particular, we keep stories hidden until they contain at least five nonduplicate articles from different sources, we hide news stories when the last update is older than eight hours, and we delete stories that stretch over a span of more than 10 days.

## Topic Classification and Entity Extraction

Since the basis of any user interaction model is a structured data model, and the source data is mostly unstructured text with few metadata items, we employ natural language processing techniques to infer additional annotations over which news items can be retrieved and visualized.

We provide two modalities of annotation: topics from a large hierarchical taxonomy, and named entities from a background knowledge base. The first structures the articles on a multilevel hierarchy of topics, and the second provides the structure across the people, organizations, and locations that appear in those articles. These two annotation mechanisms are complementary and are often employed in information management in the news domain (Sandhaus 2008), but with a manual, precision-oriented approach.

To provide these annotations, we use the Enrycher (Štajner et al. 2010) system. The topic classification is done with a hierarchical classification approach (Grobelnik and Mladenić 2004) into the DMOZ hierarchy as the data set. The entity extraction and disambiguation uses a knowledge-base approach (Štajner and Mladenić 2009) to disambiguate into DBPedia[9] (Auer et al. 2007).

## Subtopic Detection

To generate the topical map (figure 4 marker D), we need to identify the subtopics of the story and display them in a way that makes them discoverable and navigable.

To detect subtopics, we partition the articles comprising a story into at most five clusters using agglomerative k-means clustering based on bag-of-words vectors and the cosine similarity function. From each cluster, we then extract three terms with the highest tf-idf weight in its centroid and use them as the description of the subtopic. It is reasonable to assume, as we do in this approach, that each relevant subtopic will be the focus of several articles: if a subtopic is covered in equal depth in all articles, it will not form a cluster, but it is also not of particular interest for analyses of news diversity.

To map the articles (and with them, subtopics) onto a plane we use multidimensional scaling (MDS) (Fortuna, Grobelnik, and Mladenić 2005). The goal of this nonlinear dimensionality reduction technique is to position the articles in two dimensions so that their Euclidean distance in the plane agrees as well as possible with the cosine distance in the original bag-of-words space.

When the user interacts with the subtopic control, the relevance of individual articles to the user's choice is interpreted within this two-dimensional compressed topic space. Although doing so introduces some error, the result is a much simpler interaction and increased responsiveness of the controls.

## Geotagging

The map widget requires us to know the location information for individual articles. We associate each article with a news publisher, and locate those using

*Figure 5. Newsfeed.*

A monitoring system for the raw stream of articles downloaded from the Internet.

a combination of data sources in the following order of precedence:

(1) We crawled, parsed, and integrated a number of public, hand-curated news publisher listings that provide for each publisher, among other information, the city and country of origin. (2) If available, we use the country code top-level domain (ccTLD) to determine the country. For example, guardian.co.uk is mapped to the United Kingdom. (3) As a last resort, we query the WHOIS databases and heuristically parse out the address of the domain's owner. Care has to be exercised as the address has no fixed format and is besides often that of a privacy-protecting proxy registrant. In this and the previous method, we do not attempt to extract the city.

## Sentiment Detection

Sentiment analysis is a natural language processing task, which in our scenario aims to predict the polarity (positive, negative or neutral) of articles and opinions expressed therein. It is also one of the novel interaction modalities that DiversiNews offers, since the positioning of a given article's opinion on a topic is one of the big-picture questions news readers often have when confronted with new content.

To provide robust sentiment analysis, we use a supervised model combined with background knowledge in the form of sentiment lexicons. The approach (Štajner, Novalija, and Mladenić 2013) uses a multilayer classifier that first performs several independent predictions on the basis of individual feature sets (words, lexicon features, as well as orthographic features), followed by an aggregation classifier that produces the final result. We trained the classifier on an annotated news corpus (Balahur et al. 2010), combined with SentiWordNet (Baccianella, Esuli, and Sebastiani 2010), achieving an $F_1$ score of 0.78.

The classifier outputs a continuous value based on its prediction confidence to allow for better visualization and comparison of individual news articles.

## Summarization

The summaries in DiversiNews are a novel way of quickly providing users with an aggregated and tailored view of the data. We developed two multidocument summarization algorithms that can adapt to user's focus of interest.

The topic-modeling summarizer is based on replicating the latent topic distribution of the whole document collection in the summary. It is a reimplementation of the TopicSum model first described in

Haghighi and Vanderwende (2009). The algorithm aims to minimize the Kullback-Leibler (KL) divergence of a summary with respect to the document collection, that is, $KL(P_C\|P_S)$, where $P_C$ and $P_S$ are the weighted distribution of unigrams in the document collection and in the summary, respectively.

The entity-centered summarizer aims to construct the summary so as to include the most relevant entities. It extracts using subject-verb-object triplets from the documents, estimates their similarities using WordNet, and identifies the central ones with a PageRanklike algorithm. It then constructs the summary by greedily selecting sentences with the $k$ most representative triplets. Details are given in section 4 of the paper by Rusu, Trampuš, and Thalhammer (2013).

For both algorithms, we refer the reader to their respective reference papers.

### Article Ranking

When the user moves the controls in the story exploration screen, the articles are reranked to reflect the change. The ranking is computed based on the goodness of fit of each article to the position of user controls. The goodness of fit score is a linear combination of the three distances between the target value of a control and an article's value. The world map uses the logarithm of the Haversine distance, the topical map uses the Euclidean distance in the two-dimensional plane, and the sentiment slider uses the squared difference between the target sentiment and the article sentiment.

The goodness of fit is also used as an importance weight for articles when they are input into the summarizer.

# Evaluation

Most of the components have been evaluated individually before, with results presented in their respective reference papers. Here, we focus on the evaluation of the user interface as a whole and the ability of individual components to contribute to the use case of diversified news browsing.

To evaluate the reaction of users to the news browsing paradigm proposed in this article, we designed an experiment divided in two main parts: the user interface (UI) evaluation, where we measured whether the UI controls are intuitive and well-suited to the task; and, most importantly, the perceived usefulness evaluation, where the goal was to see if users find the interface useful.

The user experience evaluation was performed on the prototype web version of the interface and involved 16 test users. Each user went through the two stages of the evaluation in the same order and in the same amount of time. Of the 16 users, 14 were casual readers of web news and 2 were professionals, news operators working for a press office. All test users were volunteers. All except the two news professionals were drafted through Mechanical Turk,[10] interns at Google, or undergraduate or graduate students at the University of Karlsruhe. None of the users were members of the research groups working on the RENDER project or otherwise affiliated with the project.

Our test group roughly represents the project's target audience: For one, all the users stated they use the Internet several hours every day. On top of that, the majority (11 out of 16) of users prefer online news sites as their news channel, 2 prefer social networks, and a single test user favors newspapers.

While we extensively evaluated each component of the interface, here we limit ourselves to outlining the main findings and the lessons that we learned based on the subjects' feedback. Full results are available in a technical report (Pighin et al. 2014).

### User Interface

To assess the intuitiveness of the UI, test subjects were exposed to static images of the interface and asked to build an expectation concerning the function of the UI components without interacting with them. After that, the subjects independently used the interface for a set amount of time. Between the two activities and at the end of the second, they were asked questions about the quality of the interaction and the responsiveness and the ergonomics of the interface. This session provided useful evidence concerning the intuitiveness and accessibility of the interface.

Using dynamically changing summaries as the main form of feedback to users is unique to DiversiNews, so it was important to us to gather more user data on this aspect of the UI. To do so, we randomly selected 20 news stories and, for each of the two summarizers, we generated four different summaries based on different states of the control panels.[11] Each summary was then assessed by two evaluators. They graded fluency and informativeness to tell us something about the intrinsic quality of the summaries. In addition, they tried to identify the two out of eight summaries that correspond to current widget settings, for example, the two clearly positive-sentiment summaries. The evaluator's ability to do this correctly reflects on the topic sensitivity and sentiment sensitivity of the summarizers. We did not measure geographical sensitivity as we deemed it prohibitively hard for evaluators to identify, for example, "the Italian viewpoint" without a lot of context.

#### Controls and Integration

All the views of the interface have been found to be self-explanatory to the large majority of the subjects. The fact that the controls influence the ordering of articles and the fact that the summary synthesizes the ranked articles was clear to 75 percent or more of the users (depending on the aspect of evaluation). In table 1, question (a) shows that even before interacting with the tools, their individual intents were clear.

| | Subtopic | Publisher Location | Sentiment |
|---|---|---|---|
| (a) It was clear from the start what this control does. | 75% | 94% | 88% |
| (b) The interaction with this control is intuitive. | 69% | 81% | 81% |

*Table 1. Rater's Evalutation of Intuitiveness and Utility of Individual Focus Controls.*

The table gives percentage of users who answered "Strongly agree" or "Agree" given a five-level scale.

| | Fluency | | | Informativeness | | |
|---|---|---|---|---|---|---|
| Summarizer | Inadequate | Adequate | Human grade | Inadequate | Adequate | Human grade |
| Topic modeling | 35% | 40% | 25% | 30% | 41% | 29% |
| Entity centered | 35% | 42% | 24% | 33% | 40% | 28% |

*Table 2. Distribution of Raters' Decisions Concerning the Fluency and Informativeness of the Two Summarizers.*

The exact wording of questions and of the grading scale can be found in the technical report (Pighin et al. 2014).

Precision (P), Recall (R), and $F_1$ describe the *evaluators'* ability to identify, in a set of eight summaries, the one that the algorithm produced to be, for example, negative in sentiment.

| | Sentiment Widget | | | Topic Widget | | |
|---|---|---|---|---|---|---|
| Summarizer | *P* | *R* | $F_1$ | *P* | *R* | $F_1$ |
| Topic modeling | 0.75 | 0.59 | 0.66 | 0.98 | 0.80 | 0.88 |
| Entity centered | 0.86 | 0.53 | 0.66 | 0.75 | 0.23 | 0.35 |

*Table 3. Comparison of the Two Summarizers in Terms of Their Sensitivity to the UI Controls.*

Precision (P), Recall (R), and $F_1$ describe the *evaluators'* ability to identify, in a set of eight summaries, the one that the algorithm produced to be, for example, negative in sentiment.

Question (b) in the same table shows that interaction also proceeded as expected. The outlier in both categories is the Subtopic control. We expect this is due to two factors: first, the control navigates an inherently complex, high-dimensional aspect of data that is relatively hard to make sense of, both for humans and algorithms. Second, users are much less likely to have had prior experience with a topical map of this kind, as opposed to a geographical map or a simple slider. In this case, the complex data justifies a more complex control panel, but it is important to keep in mind that even modestly unusual interfaces can be very detrimental to ease of use and thus user adoption.

A major unexpected finding however was that during the static inspection, about half of the subjects built the expectation that acting on any of the UI controls would have an effect also on the others. For example, they imagined that changing the position of the sentiment slider would also affect the content of the topics panel to show the topics having more positive connotation. We have since altered the labeling of the widgets to further stress that the panels are independent and act on orthogonal dimensions of the data.

Summaries

The results are listed in tables 2 and 3. The two approaches are similar in performance in most metrics, with the topic modeling approach having a slight advantage. The one notable exception to this is

## How Would You Improve the Tools?
### (Aggregated over multiple focused questions)

''I have found the terms listed for the clusters not very meaningful. Some kind of summarization mechanism to generalize each cluster could help.''

''I think the mechanism for the sentiment detection should be improved.''

''The [map tool] was not always producing diverse results; otherwise fine.''

''The map is too small to use comfortably.''
*(Multiple comments in this vein.)*

''Use color to indicate the [polarity, political inclination] of elements on the [Publisher Location control, Subtopic control].'' *(Multiple comments in this vein.)*

''System could automatically generate additional dichotomies depending on topic, such as pro-life versus pro-choice in abortion debate, or pro-Samsung versus pro-Apple sentiment in patent dispute.''

''Add political/financial inclination of the newspaper news are taken from.''

''Add another control, 'uniqueness', privileging original rather than main-stream opinions.''

## What Are Your Overall Thoughts on DiversiNews?

''Excellent tool and project, but needs some UI adjustments to fly.''

''Great page, I like the summaries very much. The overall design could be improved.''

''I really like the interface and it would be really nice to use it. I would start reading more news, as the summary part seems great.''

*Figure 6. A Selection of Freeform Evaluator Feedback.*

We prompted the evaluators for suggestions for improvement or (last three answers) overall impressions.

in the summarizers' ability to adapt to topic changes; here, the topic modeling significantly outperforms the entity-centered approach. We speculate this is due to the sparse data representation used in the latter approach, which may capture the concepts associated with the main storyline, but fail to do so for the fringe concepts associated with (sub)topics.

Based on these findings, the iOS version of the news browsing interface generates summaries exclusively with the topic modeling summarizer.

## Usefulness

In this part of the evaluation, the subjects answered questions about the utility of the individual components and their potential impact on their news-browsing habits. They answered specific questions about the potential of the different components to highlight and emphasize diversity of opinion in news.

We evaluated with more than 50 questions in total and can only present an outline here; see the technical report (Pighin et al. 2014) for the full list of questions and responses. In general, the raters found that the controls succeed in modeling different dimensions and provide a more rounded paradigm for online news consumption, and found summaries to be an effective device to capture and represent relevant information and diversity of opinion. They did also suggest a number of improvements for future work.

Let us first look at our most notable areas for improvement. At multiple points during evaluation, we asked users for free-form suggestions on how the tool could be improved. At the very end, we asked for their overall thoughts. We present an aggregation of the most interesting and common responses in figure 6. There are several recurring themes: text mining errors, UI improvements, tighter integration of modalities, and new axes along which to compare news.

*Text Mining Errors.* In their sessions, the users are almost sure to have encountered articles with an incorrectly assigned subtopic or sentiment. Even if this happens for a minority of the articles, it can quickly affect users' trust in the system. We also failed to set the right expectations: many stories have little variations along one or more axes (for example, a report of a tornado will always be mostly negative in sentiment; write-ups on a tech product launch will not vary much with geography; and so on), and playing with controls in those cases will only surface algorithm or data noise. Ideally, a diversity-aware tool should be able to inform the user when data is largely homogeneous. See for example the first two responses in figure 6.

*UI improvements.* Design was not very high on our priority list, and in places this affected comfort of use. There was a lot more UI-related feedback that we do not show here, mostly regarding details. We took it into account when designing the iOS version of DiversiNews, which is now smoother to use thanks to our test users. This is important as a grating user experience can prevent users from discovering the technical worth and capabilities of a product. See for example comments three and four in figure 6.

*Tighter Integration of Modalities.* DiversiNews's controls double as rudimentary visualizations of the underlying data, but by showing for example, the sentiment superimposed on the geographical map, we could greatly improve this functionality. Users
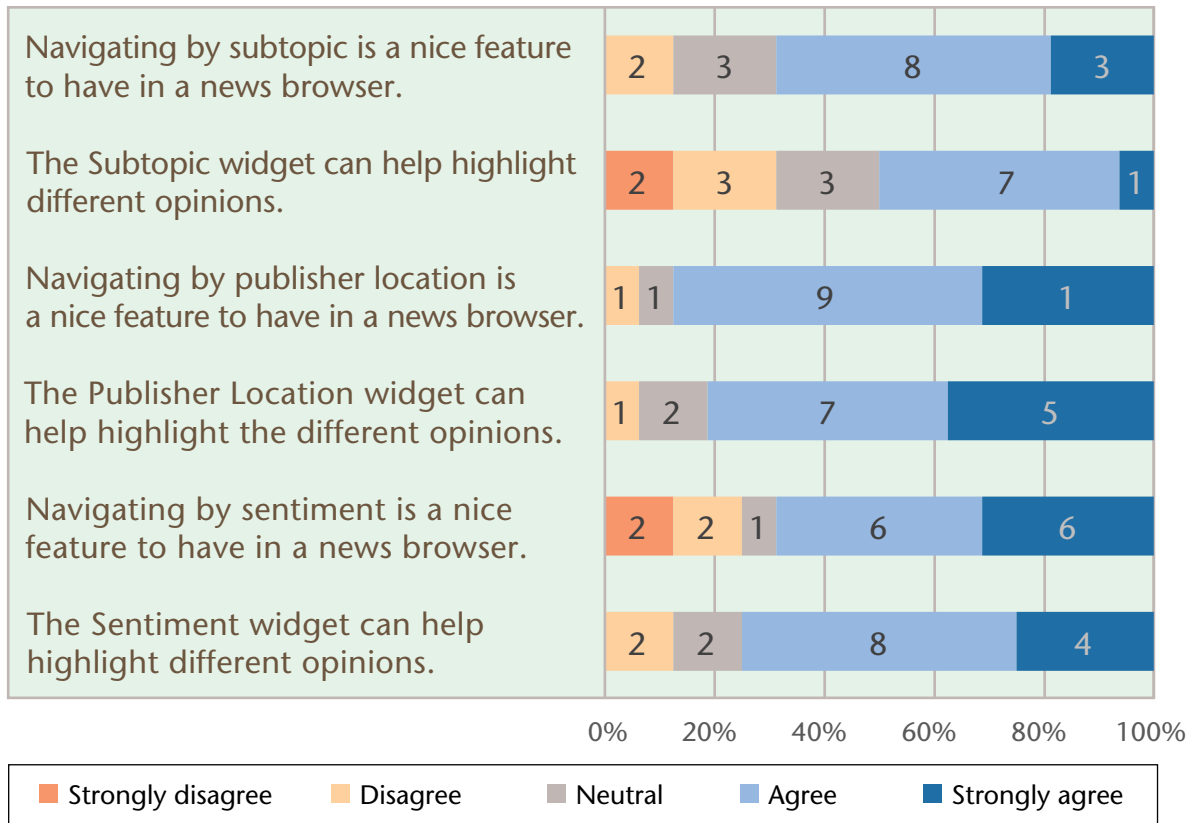
*Figure 7. Distribution of the Answers to Some of the Perceived Utility Evaluation Questions.*

Two Questions per UI Element.

mostly envisioned encoding this using color; the fifth item in figure 6 summarizes the responses in this vein.

*New Axes Along Which to Compare News.* We find this set of suggestions very encouraging, as they are completely unprompted. It is clear that there is a lot of room for further experimentation. The challenge is in implementing the new axes well, but also in not oversaturating the interface (and with that, the user). One solution to the latter is to limit the interface to a small number of the most useful dimensions. It is up to future work to determine what those are — our current set was predetermined in an internal debate as implementing and testing more controls would have been prohibitively expensive. See the final comments in figure 6 for multiple such suggestions.

Despite clear room for improvement, overall the users responded very positively to both the idea and the software. Figure 7 is central to the question of utility. Independently of our specific implementation, the majority consider browsing along each of

the axes of diversity a welcome feature for a news browser (odd rows in the table). For our specific implementation, more than 80 percent of users find that it is useful in uncovering diversity in news (even rows in the table). The Subtopic widget was received least well, both the concept and the implementation. Conceptually, we speculate the users are more interested in subjective opinion diversity rather than in nuances in focus. At the implementation level, the users found this tool the most confusing and also the most error prone in its analyses. It is far from useless however — half the users still find it useful. The widget that proved the most useful is the Source Location. Its interaction with the data is the most transparent, and its utility can be clearly seen for all politically tinged stories.

Similarly positive is the feedback to questions on summarization. More than 80 percent of the subjects found that summaries are at least adequate in quality, and just below 80 percent believe that "generating summaries according to different criteria is an effec-

tive way of letting relevant information emerge and let diversity in news stand out." The overall system comments (bottom of figure 6) are also positive — we received no strongly negative comments there.

Scope

These affirmations of usefulness need to be considered along with the scope of DiversiNews. On the one hand, when the user reads news casually, time is of value and a single short article usually fulfills the data needs quickly and conveniently. On the other hand, a system like DiversiNews requires additional focus, time, and energy, but empowers users who want or need to understand a story in depth. By stating DiversiNews is useful, our test users are not conveying "I want to read all my news with this," but rather "I have a use for this tool."

We can see DiversiNews or a similar service can be useful; but what is it that drives people to seek diverse information? Related sociology work offers some insights. At a high level, better education and higher social status correlate positively with the breadth of information people consume (Yang 2008). More specifically, Diddi and LaRose (2006) recorded news consumption patterns of college students and observed that especially people with a keen interest in current affairs tend to read more in-depth articles (a proxy for reading multiple articles as in DiversiNews). An, Quercia, and Crowcroft (2013b) aggregate related work to identify three main factors that cause diversity-seeking behavior in news consumption: (1) Weak prior beliefs and unskewed trust. If the reader has an entrenched opinion or if he trusts one source of information over others, he is less likely to seek opposing views. (2) Emotional engagement. People care more about issues that tap emotions, like worry. (3) Diverse social context. People with homogeneous social circles are more likely to get trapped in their echo chambers and not fact-check what they feel is global consensus.

## Conclusion

We developed DiversiNews, a novel news browsing system that helps users discover diverse viewpoints on a single news story, event, or issue. Users are not only presented with a large number of relevant articles from different sources, but also given the tools to manipulate them in real time according to a number of criteria. This allows them to efficiently analyze multiple perspectives in a way not possible before.

Users browse and manipulate the articles through a unified interface that succinctly presents data that is both high dimensional (subtopics) and multimodal (subtopics, geolocation, sentiment) — two properties that make data notoriously hard to navigate. The information relevant to user's expressed interests appears instantaneously and consists of a focus-aware summary, combating traditional information overload, and a list of the most pertinent arti-

cles. Both the user interface over multimodal data and the responsive, interactive textual summaries are novel to the best of our knowledge.

The solution integrates a number of both existing and purpose-made software components into a fairly complex pipeline that delivers value to the user with no special external dependencies or assumptions — starting with just an Internet connection, news articles are collected, analyzed, aggregated, and presented in a fashion that simplifies previously complex exploration that was shown to be relevant and desirable in the real world.

We evaluated the system with both casual and professional consumers of web news, receiving positive feedback from both groups. On a conceptual level, users find that making diverse news more accessible is important; on a practical level, they appreciate the summary-based interface and being in control of the criteria by which the news are organized and presented. As with any new system that introduces complexity to a common task, usability issues emerged as detailed in the Usefullness subsection; despite these problems, users see potential in DiversiNews, suggesting a future revision and reevaluation of the system are warranted.

There are still many directions in which to extend the concepts presented in this article. Carpenter (2010) and many others show that citizen journalism produces more diverse content than the mainstream media, and including blogs into DiversiNews seems promising; based on our brief experimentation, one of the challenges will be in filtering out the high amounts of chaff. It would be interesting to include time meaningfully as one of the dimensions and see how viewpoints and information change as a story evolves. As another research venue, integrating cross-lingual methods and analyses would lower the language barriers that prevent true worldwide comparison of opinions. And last but not least, much of the diversity in news can be attributed to the liberal-conservative political split, as recognized by a lot of related work; detecting the lean of articles or publishers and including it in the interface would be valuable for certain types of news stories.

## Acknowledgments

## Notes

1. Pew Research Center is a nonpartisan, nonprofit organization that conducts public opinion polling, demographic

research, media content analysis, and other empirical social science research. It is one of the most prominent U.S. organizations of its type.

2. bit.ly/1uajvHR.

3. aidemo.ijs.si/diversinews.

4. waxy.org/2008/10/memeorandum_colors.

5. emm.newsbrief.eu.

6. www.socialmention.com.

7. Mozilla Open Directory Project (dmoz.org) provides a large general-purpose hierarchy of topic categories like Sports → Soccer → Competitions → World Cup.

8. The back end serving the prototype web version is slightly simpler and uses less pre-computation and caching. However, almost all dependencies and interactions remain the same.

9. dbpedia.org .

10. Workers were selected from Google's internal roster of MTurk workers that sets a very high bar for their diligence and reliability. Workers are not affiliated with Google. Mechanical Turk can be found at mturk.com.

11. We used the extreme positions of the panels, for example, setting the Sentiment slider to maximally positive, or focusing the Subtopic panel on the center of one of the clusters.

## References

An, J.; Quercia, D.; and Crowcroft, J. 2013a. Fragmented Social Media: A Look into Selective Exposure to Political News. In *Proceedings of the 22nd international Conference on the World Wide Web*, 51–54. New York: Association for Computing Machinery.

An, J.; Quercia, D.; and Crowcroft, J. 2013b. Why Individuals Seek Diverse Opinions (Or Why They Don't). In *Web Science 2013* (WebSci '13). New York: Association for Computing Machinery. dx.doi.org/10.1145/2464464.2464493

Auer, S.; Bizer, C.; Kobilarov, G.; Lehmann, J.; Cyganiak, R.; and Ives, Z. 2007. Dbpedia: A Nucleus for a Web of Open Data. In T*he Semantic Web, 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007 + ASWC 2007*. Lecture Notes in Computer Science 4825, 722–735. Berlin: Springer. dx.doi.org/10.1007/978-3-540-76298-0_52

Azzopardi, J., and Staff, C. 2012. Incremental Clustering of News Reports. *Algorithms* 5(3): 364–378. dx.doi.org/10.3390/a5030364

Baccianella, S.; Esuli, A.; and Sebastiani, F. 2010. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2010)*. Paris: European Language Resources Association.

Balahur, A.; Steinberger, R.; Kabadjov, M. A.; Zavarella, V.; Van Der Goot, E.; Halkia, M.; Pouliquen, B.; and Belyaeva, J. 2010. Sentiment Analysis in the News. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2010)*. Paris: European Language Resources Association.

Borges, H., and Lorena, A. 2010. A Survey on Recommender Systems for News Data. In *Smart Information and Knowledge Management,* ed. E. Szczerbicki, N. T. Nguyen. Volume 260 of Studies in Computational Intelligence. Berlin: Springer.

Brzozowski, M. J.; Hogg, T.; and Szabo, G. 2008. Friends and Foes. In *Proceedings of the 2008 Conference on Human Factors in Computing Systems, CHI 2008,* 817. New York: Association for Computing Machinery. dx.doi.org/10.1145/1357054.1357183

Carpenter, S. 2010. A Study of Content Diversity in Online Citizen Journalism and Online Newspaper Articles. *New Media & Society* 12(7): 1064–1084. dx.doi.org/10.1177/1461444809348772

Diddi, A., and LaRose, R. 2006. Getting Hooked on News: Uses and Gratifications and the Formation of News Habits Among College Students in an Internet Environment. *Journal of Broadcasting & Electronic Media* 50(2): 193–210. dx.doi.org/10.1207/s15506878jobem5002_2

Downs, A. 1957. An Economic Theory of Political Action in a Democracy. *Journal of Political Economy* 65(2): 135–150. dx.doi.org/10.1086/257897

Ennals, R.; Trushkowsky, B.; and Agosta, J. 2010. Highlighting Disputed Claims on the Web. In *Proceedings of the 19th International Conference on World Wide Web*. New York: Association for Computing Machinery. dx.doi.org/10.1145/1772690.1772726

Flaxman, S.; Goel, S.; and Rao, J. 2013. Ideological Segregation and the Effects of Social Media on News Consumption. Unpublished paper available at ssrn.com/abstract=2363701.

Fortuna, B.; Grobelnik, M.; and Mladenić, D. 2005. Visualization of Text Document Corpus. *Informatica-Ljubljana* 29(4): 497.

Gamon, M.; Basu, S.; Belenko, D.; Fisher, D.; Hurst, M.; and Koenig, A. C. 2008. BLEWS: Using Blogs to Provide Context for News Articles. In *Proceedings of the Second International Conference on Weblogs and Social Media,* 60. Palo Alto, CA: AAAI Press.

Gentzkow, M., and Shapiro, J. 2011. Ideological Segregation Online and Offline. *Quarterly Journal of Economics* 126(4): 1799–1839. dx.doi.org/10.1093/qje/qjr044

Grobelnik, M., and Mladenić, D. 2004. Sim-

ple Classification into Large Topic Ontology of Web Documents. *Journal of Computing and Information Technology* 13(4): 279–285. dx.doi.org/10.2498/cit.2005.04.04

Haghighi, A., and Vanderwende, L. 2009. Exploring Content Models for Multi-Document Summarization. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*. Boulder, CO: Association for Computational Linguistics. dx.doi.org/10.3115/1620754.1620807

Kohut, A. P. R. 2013. Pew Surveys of Audience Habits Suggest Perilous Future for News. Pew Research Center Survey Report. Washington, DC: Pew Research Center.

Leetaru, K.; Wang, S.; and Cao, G. 2013. Mapping the Global Twitter Heartbeat: The Geography of Twitter. *First Monday* 18(5–6). dx.doi.org/10.5210/fm.v18i5.4366

Leskovec, J.; Kleinberg, J.; and Faloutsos, C. 2005. Graphs over Time. In *Proceeding of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining,* 177. New York: Association for Computing Machinery. dx.doi.org/10.1145/1081870.1081893

Lv, Y.; Moon, T.; Kolari, P.; Zheng, Z.; Wang, X.; and Chang, Y. 2011. Learning to Model Relatedness for News Recommendation. In *Proceedings of the 20th International Conference on World Wide Web*. New York: Association for Computing Machinery. dx.doi.org/10.1145/1963405.1963417

Maier, S. 2005. Accuracy Matters: A Cross-Market Assessment of Newspaper Error and Credibility. *Journalism & Mass Communication Quarterly* 82(3): 533–551. dx.doi.org/10.1177/107769900508200304

Milner, H. 2002. *Civic Literacy: How Informed Citizens Make Democracy Work*. Lebanon, NH: University Press of New England

Munson, S. 2012. Exposure to Political Diversity Online. Ph.D. Dissertation (Information). University of Michigan, Ann Arbor, MI.

Munson, S., and Resnick, P. 2010. Presenting Diverse Political Opinions: How and How Much. In *Proceedings of the CHI 2002 Conference on Human Factors in Computing Systems: Changing Our World, Changing Ourselves*. New York: Association for Computing Machinery.

Munson, S.; Zhou, D.; and Resnick, P. 2009. Sidelines: An Algorithm for Increasing Diversity in News and Opinion Aggregators. In *Proceedings of 3rd International AAAI Conference on Weblogs and Social Media*. Palo Alto, CA: AAAI Press.

Park, S.; Kang, S.; Chung, S.; and Song, J. 2012. A Computational Framework for Media Bias Mitigation. *ACM Transactions on*

*Interactive Intelligent Systems* 2(2): 1–32. dx.doi.org/10.1145/2209310.2209311

Park, S.; Ko, M.; Kim, J.; Choi, H.; and Song, J. 2011. NewsCube2.0: An Exploratory Design of a Social News Website for Media Bias Mitigation. Paper presented at the 2nd International Workshop on Social Recommender Systems, March 19–23, Hangzhou, China.

Park, S.; Lee, S.; and Song, J. 2010. Aspect-Level News Browsing: Understanding News Events from Multiple Viewpoints. In *Proceedings of the 2010 International Conference on Intelligent User Interfaces.* New York: Association for Computing Machinery. dx.doi.org/10.1145/1719970.1719977

Pighin, D.; Alfonseca, E.; Leif Keppmann, F.; and Trampuš, M. 2014. Evaluation of the DiversiNews Diversified News Service. Technical report arXIV:1407.4454. Ithaca, NY: Cornell University Library.

Rortais, A.; Belyaeva, J.; Gemo, M.; Goot, E. V. D.; and Linge, J. P. 2010. MedISys: An Early-Warning System for the Detection of (Re-) Emerging Food- and Feed-Borne Hazards. *Food Research International* 43(5): 1553–1556. dx.doi.org/10.1016/j.foodres.2010. 04.009

Rusu, D.; Trampuš, M.; and Thalhammer, A. 2013. Diversity-Aware Summarization — RENDER Project Deliverable D3.2.1. Technical Report RENDER project. Karlsruhe, Germany: Karlsruhe Institute of Technology.

Sandhaus, E. 2008. *The New York Times Annotated Corpus.* Linguistic Data Consortium, Philadelphia 6(12): e26752.

Silva, J. A.; Faria, E. R.; Barros, R. C.; Hruschka, E. R.; de Carvalho, A. C. P. L. F.; and Gama, J. a. 2013. Data Stream Clustering: A Survey. *ACM Computing Survey* 46(1): 13:1–13:31.

Štajner, T., and Mladenić, D. 2009. Entity Resolution in Texts Using Statistical Learning and Ontologies. In *The Semantic Web: Proceedings of the 8th International Semantic Web Conference, ISWC 2009*. Lecture Notes in Computer Science Volume 5823, 91–104. Berlin: Springer. dx.doi.org/10.1007/978-3-642-10871-6_7

Štajner, T.; Novalija, I.; and Mladenić, D. 2013. Informal Multilingual Multi-Domain Sentiment Analysis. *Informatica* 37(4): 373–380.

Stajner, T.; Rusu, D.; Dali, L.; Fortuna, B.; Mladenić, D.; and Grobelnik, M. 2010. A Service Oriented Framework for Natural Language Text Enrichment. *Informatica* 34(3): 307–313.

Steinberger, R.; Pouliquen, B.; and Goot, E. V. D. 2009. An Introduction to the Europe Media Monitor Family of Applications. Paper presented at the SIGIR 2009 Workshop on Information Access in a Multilingual World, Boston, MA July 23.

Tavakolifard, M.; Gulla, J. A.; Almeroth, K. C.; Ingvaldsn, J. E.; Nygreen, G.; and Berg, E. 2013. Tailored News in the Palm of Your Hand: A Multiperspective Transparent Approach to News Recommendation. In *Proceedings of the 22nd International Conference on the World Wide Web,* 305–308. New York: Association for Computing Machinery.

Trampuš, M., and Novak, B. 2012. Internals of an Aggregated Web News Feed. Paper presented at the Fifteenth International Multiconference on Information Society, Ljubljana, Slovenia, 8–12 October.

Voakes, P., and Kapfer, J. 1996. Diversity in the News: A Conceptual and Methodological Framework. *Journalism & Mass Communication Quarterly* 73(3): 582–593. dx.doi.org/10.1177/107769909607300306

Yang, J. 2008. The Widening Information Gap Between High and Low Education Groups: Knowledge Acquisition from Online Versus Print News. Doctoral diss., Department of Mass Communications and Telecommunications, Indiana University Bloomington, Bloomington, IN.

Zhang, J.; Kawai, Y.; and Kumamoto, T. 2012. Extracting Similar and Opposite News Websites Based on Sentiment Analysis. Paper presented at the 2012 International Conference on Industrial and Intelligent Information (ICIII 2012), Singapore, 17–18 March.

**Mitja Trampuš** is a Ph.D. candidate in text mining. He majored in computer science and mathematics at University of Ljubljana in 2008 and has since worked on several international research projects at Jozef Stefan Institute. He gained further experience during internships with data science teams at Facebook, Twitter, and Google and currently works at Twitter.

**Flavio Fuart** spent the last nine years developing software for online media monitoring and open-source intelligence dealing with productization, optimization, integration, and deployment of research modules developed by his colleagues. Occasionally he carries on research tasks as well. He holds a university degree in computer science from the Faculty of Computer and Information Science, University of Ljubljana, Slovenia (1997). During his career, he worked for SMEs (1994–2005), European Commission (2005–2012), and United Nations (2009), and currently he works at Josef Stefan Institute, Ljubljana, Slovenia.

**Daniele Pighin** is a research scientist at Google Zurich, where his main contributions are in the areas of relation extraction and event learning. Before moving to Switzerland, he worked on machine translation and its evaluation, kernel methods for natural language processing, and semantic role labeling.

**Tadej Štajner** is a Ph.D. student in artificial intelligence with a focus on natural language processing. His research applies to multilingual and cross-lingual text analysis with the goal of better understanding how to automatically learn text representations that allow for a language-independent usage in applications. He is also a member of the W3C MultilingualWeb-LT Working Group and a coauthor of the Internationalization Tag Set 2.0 standard.

**Jan Berčič** is a bachelor in mathematics and is currently studying computer science at the University of Ljubljana. He has some experience with competitive programming and systems engineering and is working for the Artificial Intelligence Laboratory of the Jozef Stefan Institute.

**Blaz Novak** has received his BSc in computer science from University of Ljubljana in 2008. He is currently working as a researcher at Jozef Stefan Institute, pursuing a Ph.D. in machine learning and data mining.

**Delia Rusu** is a doctoral student at the Jozef Stefan International Postgraduate School. Her dissertation topic is the automatic annotation of text with background knowledge, using statistical models of natural language and semantic technologies.

**Luka Stopar** graduated at the Faculty of Computer and Information Science at the University of Ljubljana (ISRM). He is a research assistant at Jozef Stefan Institute. He is working toward a Ph.D. in text mining at the Jozef Stefan International Postgraduate School.

**Marko Grobelnik** is an expert in the analysis of large amounts of complex data with the purpose to extract useful knowledge. Apart from research on theoretical aspects of unconventional data analysis techniques, he has valuable experience in the field of practical applications and development of business solutions based on innovative technologies. Marko works with industry leaders such as Microsoft Research, the *New York Times*, and British Telecom among others. He has published papers regularly in refereed conferences and journals.