

# An Infrastructure for Empowering Internet Users to Handle Fake News and Other Online Media Phenomena

Georg Rehm<sup>(✉)</sup>

DFKI GmbH, Language Technology Lab,  
Alt-Moabit 91c, 10559 Berlin, Germany  
[georg.rehm@dfki.de](mailto:georg.rehm@dfki.de)

**Abstract.** Online media and digital communication technologies have an unprecedented, even increasing level of social, political and also economic relevance. This article proposes an infrastructure to address phenomena of modern online media production, circulation and manipulation by establishing a distributed architecture for automatic processing and human feedback.

## 1 Introduction

The umbrella term “fake news” is often used to refer to a number of different phenomena around online media production, circulation, reception and manipulation that emerged in recent years and that have been receiving a lot of attention from multiple stakeholders including politicians, journalists, researchers, non-governmental organisations, industry and civil society. In addition to the challenge of dealing with “fake news”, “alternative facts” as well as “post-truth politics”, there is an increasing amount of hate speech, abusive language and cyber bullying taking place online.

Among the interested stakeholders are politicians who have begun to realise that, increasingly, major parts of public debates and social discourse are carried out online, on a small number of social networks. We have witnessed that not only online discussions but also the perception of trends, ideas, theories, political parties, individual politicians, elections and societal challenges can be subtly influenced and significantly rigged using targeted social media campaigns, devised at manipulating opinions to create long-term sustainable mindsets on the side of the recipients. We live in a time in which online media, online news and online communication have an unprecedented level of social, political and economic relevance.

Due to the intrinsic danger of successful large-scale manipulations the topic is of utmost importance. Many researchers from the Social Sciences and Computer Science currently work on the topic. An idea often mentioned is to design, develop and deploy technologies to improve the situation, maybe even to solve it altogether, thanks to recent breakthroughs in AI (Metz 2016; Gershgorn 2016;

Martinez-Alvarez 2017; Chan 2017), while at the same time *not* putting in place a centralised infrastructure, which could be misused for censorship, manipulation or mass surveillance.<sup>1</sup>

This article addresses key challenges of the digital age (Sect. 2) by introducing and proposing the vision of a technological infrastructure (Sect. 3); the concept has been devised in a research and technology transfer project, in which smart technologies for curating large amounts of digital content are being developed and applied by companies that cover different sectors including journalism (Rehm and Sasaki 2015; Bourgonje et al. 2016a,b; Rehm et al. 2017). Among others, we currently develop services aimed at the detection and classification of abusive language (Bourgonje et al. 2017a) and clickbait content (Bourgonje et al. 2017b). The proposed hybrid infrastructure combines automatic language technology components and user-generated annotations and is meant to empower internet users better to handle the modern online media phenomena mentioned above.

## 2 Modern Online Media Phenomena

The World Wide Web makes it possible for everybody to create content, to write an article on a certain topic. Until a few years ago the key challenge was to optimise the HTML code, linking and metadata to get highly ranked by the relevant search engines. Nowadays content is no longer predominantly discovered through search engines but through social media platforms: users see interesting content, which is then shared to their own connections. Many users only read a headline, identify a certain relevance to their own lives and then share the content. When in doubt, users estimate the trustworthiness of the source: potentially dubious stories about which they are skeptical are shared anyway if the friend through whom the story was discovered is considered reliable or if the number of views is rather high, which, to many users, indicates legitimacy.

There is a tendency for provocative, aggressive, one-sided, allegedly “authentic” (Marchi 2012) content. The idea is to make it as easy as possible to identify the stance of the article so that the reader’s own world view is validated, implicitly urging the user to share the content. The publisher’s goal is for a story to go viral, that it is shared rapidly by many users and spread through the networks to establish a reach of millions. One sub-category of this type of content is “clickbait”, articles with dubious factual content, presented with misleading headlines, designed for the simple goal of generating many views. The more extreme the virality, the higher the reach, the higher the click numbers, the higher the advertisement revenue. The term “clickbait” can also refer to articles spreading political mis- or disinformation.

Content is typically discovered through a small number of social networks. While search engines and online news portals used to be the central points of

<sup>1</sup> An indicator for the relevance of the topic is the increasing number of “how to identify fake news” articles published online (Mantzaris 2015; Bazzaz 2016; Rogers and Bromwich 2016; Wardle 2017; Walbrühl 2017).

information until a few years ago, the role of the centralised hub – and gatekeeper – is now played by social networks that help content to be discovered and go viral (Barthel et al. 2016). All social networks have as their key feature a news feed or timeline, i.e., posts, news, ads, tweets, photos presented to the user, starting with the most recent one. Nearly all social networks use machine learning algorithms to determine which content to present to a certain user. They are continuously trained through interactions with the network, i.e., “liking” a post boosts the respective topic, visiting the profile of a “friend” boosts this connection. Some networks use more fine-grained sentiments in addition to the simple “like” (see, e.g., Facebook’s reactions “love”, “haha”, “wow”, “sad”, “angry”). Through “likes” of topics, connections to friends and interactions with the site, social networks create, and continuously update, for every single user, a model of their interests, which is used to select content for the user’s timeline. The algorithms are designed to favour content liked or shared by those friends the user interacts with the most. This is the origin of the filter bubble phenomenon: users are predominantly exposed to content that can also be described as “safe” and “non-controversial” – content shared by friends they know and like is considered content that matches a user’s interests. Content that contradicts a user’s world view or that challenges their beliefs is *not* presented.

Additionally, we are faced with the challenge that more and more content is produced and spread with the sole purpose of manipulating the readers’ beliefs and opinions by appealing to their emotions instead of informing them objectively. Rather, this type of opinionated, emotional, biased, often aggressive and far-right content is spread to accomplish specific goals, for example, to create support for controversial ideas or to intensify the division between two social groups. These coordinated campaigns are carried out by experts with in-depth knowledge of the underlying technologies and processes. They involve large numbers of bots and fake accounts as amplifiers (Weedon et al. 2017) as well as large budgets for online advertisements in social media, clearly targeted at very specific demographic groups the originators want to influence and then to flip to reach a specific statistical threshold. The way news are nowadays spread, circulated, consumed and shared – with less and less critical thinking or fact checking – enables this type of content to gather a large number of readers (and sharers) quickly. The filter bubble acts like an echo chamber that can amplify any type of content, from genuine, factual news to emotionally charged, politically biased news, to false news to orchestrated disinformation campaigns, created with the specific purpose of large-scale manipulation. Content of the last two categories can be hard or very hard to identify even for human experts.

A key challenge for users and machines alike is to separate objective, balanced content, be it journalistic or user-generated, from hateful, abusive or biased content, maybe produced with a hidden agenda. Even if fundamentally different in nature, both types of content share the same potential level of visibility, reach and exposure through the equalisation mechanisms of the social web, which is prone to manipulation. In the past the prerequisite tasks of fact checking, critical thinking and uncovering hidden agendas have been in the realm of (investigative)

journalism – in the digital age they are more and more transferred to the actual reader of online content. The analysis and assessment of content is no longer carried out by professional journalists or news editors – the burden of fact checking and content verification is left to the reader. This aspect is getting even more crucial because the number of people who state that social networks are their *only* source of news is growing steadily (Marchi 2012). The most prominent example from recent history is the ongoing debate whether highly targeted social media ads influenced the 2016 US presidential election (Barthel et al. 2016; Rogers and Bromwich 2016; Marwick and Lewis 2017). It must be noted that a large number of fact-checking initiatives is active all over the world (Mantzaris 2017) but they mostly rely on human expertise and, thus, do not scale (Martinez-Alvarez 2017; Dale 2017). The small number of automated fact checking initiatives are fragmented (Babakar and Moy 2016).

**Table 1.** Characteristics and intentions associated with different types of false news – adapted from (Wardle 2017; Walbrühl 2017; Rubin et al. 2015; Holan 2016; Weedon et al. 2017)

	Satire or parody	False connection	Misleading content	False context	Imposter content	Manipulated content	Fabricated content
Clickbait		X	X	?		?	?
Disinformation			X	X		X	X
Politically biased		?	X	?		?	X
Poor journalism		X	X	X			
To parody	X				?		X
To provoke					X	X	X
To profit	?	X			X		X
To deceive		X	X	X	X	X	X
To influence politics			X	X		X	X
To influence opinions			X	X	X	X	X

Several types of online content are often grouped together under the label “fake news”. For example, Holan (2016) defines fake news as “made-up stuff, masterfully manipulated to look like credible journalistic reports that are easily spread online to large audiences willing to believe the fictions and spread the word.” In reality, the situation is much more complex. Initially based on the classification suggested by Wardle (2017), Table 1 shows an attempt at bringing together the different types of false news including selected characteristics and associated intentions. The table shows the complexity of the situation and that a more fine-grained terminology is needed to discuss the topic properly, especially when it comes to designing technological solutions that are meant to address one or more of these types of content.

An additional challenge is the proliferation of hateful comments and abusive language, often used in the comments and feedback sections on social media posts. The effects can be devastating for the affected individual. Many hateful

comments on repeated postings by the same person, say, a pupil, are akin to cyberbullying and cybermobbing. There is also a clear tendency to aggressive comments on, for example, the social media pages of traditional news outlets, who have to ask the users more and more to behave in a civilised way.

### 3 Technology Framework: Approach

Technically, online content is predominantly consumed through two possible channels, both of which rely substantially on World Wide Web technology and established web standards. Users either read and interact with content directly on the web (mobile or desktop versions of websites) or through dedicated mobile apps; this can be considered using the web implicitly as many apps make heavy use of HTML5 and other web technologies. The World Wide Web itself still is and, for the foreseeable future, will continue to be the main transport medium for online content. The infrastructure suggested by this article is, hence, designed as an additional layer on top of the World Wide Web. The scope and ambition of the challenge is immense because the infrastructure needs to be able to cope with millions of users, arbitrary content types, hundreds of languages and massive amounts of data. Its goal is to empower users by enabling them to balance out network and filter bubble effects and to provide mechanisms to filter for abusive content.

#### 3.1 Services of the Infrastructure

The burden of analysing and fact checking online content is often shifted to the reader (Sect. 2), which is why corresponding analysis and curation services need to be made available in an efficient and ubiquitous way. The same tools to be used by *content consumers* can and should also be applied by *content creators*, e.g., journalists and bloggers. Those readers who are interested to know more about what they are currently reading should be able to get the additional information as easily as possible, the same applies to those journalists who are interested in fact-checking the content they are researching for the production of new content.

Readers of online content are users of the World Wide Web. They need, first and foremost, web-based tools and services with which they can process any type of content to get additional information on a specific piece, be it one small comment on a page, the main content component of a page (for example, an article) or even a set of interconnected pages (one article spread over multiple pages), for which an assessment is sought.

The services need to be designed to operate in and with the web stack of technologies, they need to support users in their task of reading and curating content within the browser in a smarter and, eventually, more balanced way. This can be accomplished by providing additional, also alternative opinions and view points, by presenting other, independent assessments, or by indicating if content is dangerous, abusive, factual or problematic in any way. Fully automatic technologies (Rubin et al. 2015; Schmidt and Wiegand 2017; Horne and Adal

2017; Martinez-Alvarez 2017) can take over a subset of these tasks but, given the current state of the art, not all, which is why the approach needs to be based both on simple and complex automatic filters and watchdogs as well as human intelligence and feedback.<sup>2</sup>

The tools and services should be available to every web user without the need of installing any additional third-party software. This is why the services, ideally, should be integrated into the browser on the same level as bookmarks, the URL field or the navigation bar, i.e., without relying on the installation of a plugin. The curation services should be thought of as an inherent technology component of the World Wide Web, for which intuitive and globally acknowledged user-interface conventions can be established, such as, for example, traffic light indicators for false news content (green: no issues found; yellow: medium issues found and referenced; red: very likely false news). Table 2 shows a first list of tools and services that could be embedded into such a system.<sup>3</sup> Some of these can be conceptualised and implemented as automatic tools (Horne and Adal 2017), while others need a hybrid approach that involves crowd-sourced data and opinions. In addition to displaying the output of these services, the browser interface needs to be able to gather, from the user, comments, feedback, opinions and sentiments on the current piece of content, further to feed the crowd-sourced data set. The user-generated data includes both user-generated annotations (UGA) and also user-generated metadata (UGM). Automatically generated metadata are considered machine-generated metadata (MGM).

**Table 2.** Selected tools and services to be provided through the infrastructure

Tool or service	Description	Approach
Political bias indicator	Indicates the political bias (Martinez-Alvarez 2017) of a piece of content, e.g., from far left to far right	Automatic
Hate speech indicator	Indicates the level of hate speech a certain piece of content contains	Automatic
Reputation indicator	Indicates the reputation, credibility (Martinez-Alvarez 2017), trustworthiness, quality (Filloux 2017) of a certain news outlet or individual author of content	Crowd, automatic
Fact checker	Checks if claims are backed up by references, evidence, established scientific results and links claims to the respective evidence (Babakar and Moy 2016)	Automatic
Fake news indicator	Indicates if a piece of content contains non-factual statements or dubious claims (Horne and Adal 2017; Martinez-Alvarez 2017)	Crowd, automatic
Opinion inspector	Inspect opinions and sentiments that other users have with regard to this content (or topic) – not just the users commenting on one specific site but all of them	Crowd, automatic

<sup>2</sup> A fully automatic solution would work only for a very limited set of cases. A purely human-based solution would work but required large amounts of experts and, hence, would not scale. This is why we favour, for now, a hybrid solution.

<sup>3</sup> This list is meant to be indicative rather than complete. For example, services for getting background information on images are not included (Gupta et al. 2013). Such tools could help pointing out image manipulations or that an old image was used, out of context, to illustrate a new piece of news.

### 3.2 Characteristics of the Infrastructure

In order for these tools and services to work effectively, efficiently and reliably, they need to have several key characteristics, which are critical for the success of the approach.

Like the Internet and the World Wide Web, the proposed infrastructure must be operated in a federated, i.e., de-centralised setup – a centralised approach would be too vulnerable for attacks or misuse. Any organisation, company, research centre or NGO should be able to set up, operate and offer services (Sect. 3.1) and pieces of the infrastructure. The internal design of the algorithms and tools may differ but their output should comply to a standardised metadata format (MGM). It is rather likely that political biases in different processing models meant to serve the same purpose cannot be avoided, which is especially likely for models based on large amounts of data, which, in turn, may inherently include a political bias. This is why users must be enabled to configure their own personalised set of tools and services to get an aggregated value, for example, with regard to the level of hate speech in content or its political bias. Services and tools must be combinable, i.e., they need to comply to standardised input and output formats (Babakar and Moy 2016). They also need to be transparent (Martinez-Alvarez 2017). Only transparent, i.e., fully documented, checked, ideally also audited approaches can be trustworthy.

Access to the infrastructure should be universal and available everywhere, i.e., in any browser, which essentially means that, ideally, the infrastructure should be embedded into the technical architecture of the World Wide Web. As a consequence, access mechanisms should be available in every browser, on every platform, as native elements of the GUI. These functions should be designed in such a way that they support users without distracting them from the content. Only if the tools are available virtually anywhere, can the required scale be reached.

The user should be able to configure and to combine multiple services, operated in a de-centralised way, for a clearly defined purpose in order to get an aggregated value. There is a danger that this approach could result in a replication and shift of the filter bubble effect (Sect. 2) onto a different level but users would at least be empowered actively to configure their own personal set of filters to escape from any resulting bubble. The same transparency criterion also applies to the algorithm that aggregates multiple values.

### 3.3 Building Blocks of the Proposed Infrastructure

Research in Language Technology and NLP currently concentrates on smaller components, especially watchdogs, filters and classifiers (see Sect. 4) that could be applied under the umbrella of a larger architecture to tackle current online media phenomena (Sect. 2). While this research is both important and crucial, even if fragmented and somewhat constrained by the respective training data sets (Rubin et al. 2015; Conroy et al. 2015; Schmidt and Wiegand 2017) and limited use cases, we also need to come to a shared understanding *how* such

components can be deployed and made available. The proposed infrastructure consists of several building blocks (see Fig. 1).

**Building Block: Natively embedded into the World Wide Web** – An approach that is able to address modern online media and communication phenomena adequately needs to operate on a web-scale level. It should natively support cross-lingual processing and be technically and conceptually embedded into the architecture of the World Wide Web itself. It should be standardised, endorsed and supported not only by all browser vendors but also by all content and media providers, especially the big social networks and content hubs. Only if *all* users have *immediate* access to the tools and services suggested in this proposal can they reach its full potential. The services must be unobtrusive and cooperative, possess intuitive usability, their recommendations and warnings must be immediately understandable, it must be simple to provide general feedback (UGM) and assessments on specific pieces of content (UGA).

**Building Block: Web Annotations** – Several pieces of the proposed infrastructure are already in place. One key component are Web Annotations, standardised by the World Wide Web Consortium (W3C) in early 2017 (Sanderson et al. 2017a,b; Sanderson 2017). They enable users to annotate arbitrary pieces of web content, essentially creating an additional and independent layer on top of the regular web. Already now Web Annotations are used for multiple individual projects in research, education, scholarly publishing, administration and investigative journalism.<sup>4</sup> Web Annotations are *the* natural mechanism to enable users and readers interactively to work with content, to include feedback and assessments, to ask the author or their peers for references or to provide criticism. The natural language content of Web Annotations (UGA) can be automatically mined using methods such as sentiment analysis or opinion mining – in order to accomplish this across multiple languages, cross-lingual methods need to be applied (Rehm et al. 2016). However, there are still limitations. Content providers need to enable Web Annotations by referencing a corresponding JavaScript library. Federated sets of annotation stores or repositories are not yet foreseen, neither are native controls in the browser that provide aggregated feedback, based on automatic (MGM) or manual content assessments (UGM, UGA). Another barrier for the widespread use and adoption of Web Annotations are proprietary commenting systems, as used by all major social networks. Nevertheless, services such as Hypothes.is enable Web Annotations on any web page, but native browser support, ideally across all platforms, is still lacking. A corresponding browser feature needs to enable both free-text annotations of arbitrary content pieces (UGA) but also very simple flagging of problematic content, for example, “content pretends to be factual but is of dubious quality” (UGM). Multiple UGA, UGM or MGM annotations could be aggregated and presented to new readers of the content to provide guidance and indicate issues.

---

<sup>4</sup> See, for example, the projects presented at I Annotate 2015 (<http://iannotate.org/2015/>), 2016 (<http://iannotate.org/2016/>) and 2017 (<http://iannotate.org/2017/>).



**Building Block: Metadata Standards** – Another needed piece of the architecture is an agreed upon metadata schema, i.e., a controlled vocabulary, (Babakar and Moy 2016) to be used both in manual annotation scenarios (UGM) and also by automatic tools (MGM). Its complexity should be as low as possible so that key characteristics of a piece of content can be adequately captured and described by humans or machines. With regard to this requirement, W3C published several standards to represent the provenance of digital objects (Groth and Moreau 2013; Belhajjame et al. 2013a). These can be thought of as descriptions of the entities or activities involved in producing or delivering a piece of content to understand how data was collected, to determine ownership and rights or to make judgements about information to determine whether to trust content (Belhajjame et al. 2013b). An alternative approach is for publishers to use Schema.org’s ClaimReview<sup>5</sup> markup after specific facts have been checked. The needed metadata schema can be based on the W3C provenance ontology and/or Schema.org. Additional metadata fields are likely to be needed.

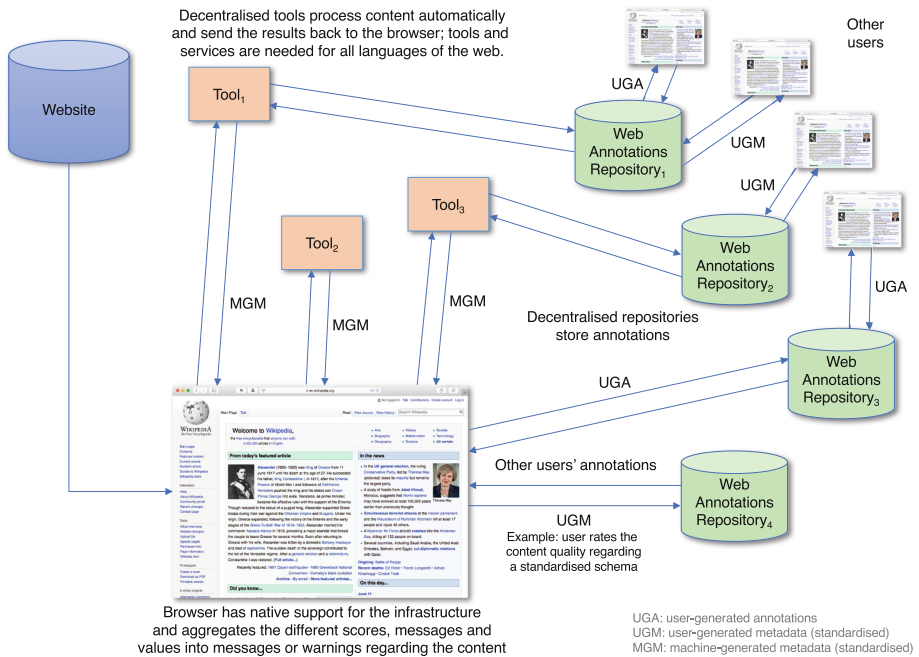
**Building Block: Tools and Services** – Web Annotations can be used by readers of online content to provide comments or to include the results of researched facts (UGA, UGM). Automatic tools and services that act as filters or watchdogs can make use of the same mechanisms (MGM, see Sect. 3.1). These could be functionally limited classifiers, for example, regarding abusive language, or sophisticated NLU components that attempt to check certain statements against one or more knowledge graphs. Regardless of the complexity and approach, the results can be made available as globally accessible Web Annotations (that can even, in turn, be annotated themselves). Services and tools need to operate in a decentralised way, i.e., users must be able to choose from a wide variety of automatic helpers. These could, for example, support users to position content on the political spectrum, either based on crowd-sourced annotations, automatic tools, or both (see Table 2).

**Building Block: Decentralised Repositories and Tools** – The setup of the infrastructure must be federated and decentralised to prevent abuse by political or industrial forces. Data, especially annotations, must be stored in decentral repositories, from which browsers retrieve, through secure connections, data to be aggregated and displayed (UGM, UGA, MGM, i.e., annotations, opinions, automatic processing results etc.). In the medium to long term, in addition to annotations, repositories will also include more complex data, information and knowledge that tools and services will make use of, for example, for fact checking. In parallel to the initiative introduced in this article, crowd-sourced knowledge graphs such as Wikidata or DBpedia will continue to grow, the same is true for semantic databases such as BabelNet and many other data sets, usually available and linkable as Linked Open Data. Already now we can foresee more sophisticated methods of validating and fact-checking arbitrary pieces of content using systems that make heavy use of knowledge graphs, for example, through automatic entity recognition and linking, relation extraction, event extraction and

<sup>5</sup> <https://schema.org/ClaimReview>.

mapping etc. One of the key knowledge bases missing, in that regard, is a Web Annotation-friendly event-centric knowledge graph, against which fact-checking algorithms can operate.<sup>6</sup> Basing algorithms that are supposed to determine the truth of an arbitrary statement on automatically extracted and formally represented knowledge creates both practical and philosophical questions, among others, who checks these automatically extracted knowledge structures for correctness? How do we represent conflicting view points and how do algorithms handle conflicting view points when determining the validity of a statement? How do we keep the balance between multiple subjective opinions and an objective and scientific ground-truth?

**Building Block: Aggregation of Annotations** – The final key building block of the proposed infrastructure relates to the aggregation of automatic and manual annotations, created in a de-centralised and highly distributed way by human users and automatic services (UGA, UGM, MGM). Already now we can foresee very large numbers of annotations so that the aggregation and consolidation will be a non-trivial technical challenge. This is also true for those human annotations that are not based on shared metadata vocabularies but that are free text – for these free and flexible annotations, robust and also multilingual annotation mining methods need to be developed.



**Fig. 1.** Simplified architecture of the proposed infrastructure

<sup>6</sup> Promising candidates are GDELT (<http://www.gdelproject.org>) and EventRegistry (<http://eventregistry.org>).

## 4 Related Work

Research on Computer-Mediated Communication (CMC) has a long tradition. Scholars initially concentrated on different types of communication media such as e-mail, IRC, Usenet newsgroups, and different hypertext systems and document types, especially personal home pages, guestbooks and, later, discussion fora (Runkehl et al. 1998; Crystal 2001; Rehm 2002). Early on, researchers focused upon the (obvious) differences between these new forms of digital communication and traditional forms, especially when it comes to linguistic phenomena that can be observed on the text surface (smileys, emoticons, acronyms etc.). Several authors pointed out that the different forms of CMC have a certain oral and spoken style, quality and conceptualisation to them, as if produced spontaneously in a casual conversation, while being realised in a written medium (Haase et al. 1997).

If we now fast forward to 2017, a vastly different picture emerges. About half of the global population has access to the internet, most of whom also use the World Wide Web and big social networks. Nowadays the internet acts like an amplifier and enabler of social trends. It continues to penetrate and to disrupt our lives and social structures, especially our traditions of social and political debates. The relevance of online media, online news and online communication could not be any more crucial. While early analyses of CMC, e.g., (Reid 1991), observed that the participants were involved in the “deconstruction of boundaries” and the “construction of social communities”, today the exact opposite seems to be the case: not only online but also offline can we observe the trend of increased, intricately orchestrated, social and political manipulation, nationalism and the exclusion of foreigners, immigrants and seemingly arbitrary minorities – boundaries are constructed, social communities deconstructed, people are manipulated, individuals excluded.

There is a vast body of research on the processing of online content including text analytics (sentiment analysis, opinion and argument mining), information access (summarisation, machine translation) and document filtering (spam classification), see (Dale 2017). Attempting to classify, among others, the different types of false news shown in Table 1 requires, as several researchers also emphasise, a multi-faceted approach that includes multiple different processing steps. We have to be aware of the ambition, though, as some of the “fake news detection” use case scenarios are better described as “propaganda detection”, “disinformation detection”, maybe also “satire detection”. These are difficult tasks at which even humans often fail. Current research in this area is still fragmented and concentrates on very specific sub-problems, see, for example, the Fake News Challenge, the Abusive Language Workshop, or the Clickbait Challenge.<sup>7</sup> What is missing, however, is a practical umbrella that pulls the different pieces and resulting technology components together and that provides an approach that can be realistically implemented and deployed including automatic tools as well as human annotations.

<sup>7</sup> See <http://www.fakenewschallenge.org>, <http://www.clickbait-challenge.org>, <https://sites.google.com/site/abusivelanguageworkshop2017/>.

## 5 Summary and Conclusions

Humanity is transitioning into becoming a digital society, or at least a “digital first” society, i.e., news, media, facts, rumours (Zubiaga et al. 2016; Derczynski et al. 2017; Srivastava et al. 2017), information are created, circulated and consumed online. Already now the right social media strategy can make or break an election or influence if a smaller or larger societal or demographic group (city, region, country, continent) is in favour or against constructively solving a certain societal challenge. Social media and online communication technologies can be an extremely powerful tool to bridge barriers, inform people and enable global communication and a constructive dialogue. When abused, misused or infiltrated, they are a dangerous weapon.

Computational Linguistics, Language Technology and Artificial Intelligence should actively contribute solutions to this key challenge of the digital age. If not, there is a concrete danger that stakeholders with bad intentions are able to influence parts of the society to their liking, only constrained by their political or commercial interests. Technologies need to be developed to enable every user of online media to break out of their filter bubbles and to inform themselves in a balanced way, taking all perspectives into account. Nevertheless, there is, as Dale (2017) points out, the danger that technologies developed to *detect* false news can also be used to *create* false news.

After dumb digital content, smart content and semantic content enrichment we now need to concentrate on content curation tools that enable *contextualised content*, i.e., content that can be, ideally, automatically cross-referenced and fact-checked, and for which background information can be retrieved in a robust way. This can involve assessing the validity of claims as well as retrieving related texts, facts and statements, both in favour and against a certain piece of content.

In this article a hybrid technology infrastructure that provides user- and machine-generated annotations on top of the whole World Wide Web is proposed with the ultimate goal of empowering internet users to handle false news and other online media phenomena by providing both automatic assessments of content and also by including alternative opinions into the process of media consumption. However, part of the solution could be provided by the small number of social networks which currently connect a vast majority of the online population and whose features and mechanisms are responsible for and also amplify the phenomena discussed in this article. It can be argued that these social networks have an obligation to act, for example, by modifying their algorithms to enable users to break out of their filter bubbles, by making the algorithms more transparent, or by using data analytics to detect potential manipulations. It is likely that regulatory steps will be taken on national and international levels.

Future work includes presenting this proposal in various different fora and communities, among others, researchers and technologists, standards-developing organisations (Babakar and Moy 2016) and national as well as international political bodies. At the same time, research needs to be continued and prototypes of the architecture as well as individual services developed, enabling organisations to build and to deploy decentralised tools early. While a universal, globally

accessible, balanced and well maintained knowledge graph containing up-to-date information about entities and events would be convenient to have, it is out of scope with regard to the initiative proposed in this article; however, it is safe to assume that such a knowledge repository will be developed in parallel in the next couple of years. The proposed infrastructure can be used to link online content against such a knowledge graph and to measure the directions of online debates.

The proposal introduced in this article is ambitious in its scope and implications, prevention of misuse and establishing trust will play a hugely important role. How can we make sure that a certain piece of technology is only used with good intentions? Recently it has been shown that a user's social media data can reliably predict if the user is suffering from alcohol or drug abuse (Ding et al. 2017). Will this technology be used to help people or to stigmatise them? Will an infrastructure, as briefly sketched in this paper, be used to empower users to make up their own minds by providing additional information about online content or will it be used to spy on them and to manipulate them with commercial or political intentions?

**Acknowledgments.** The author would like to thank the reviewers for their insightful comments and suggestions. The project “Digitale Kuratierungstechnologien” (DKT) is supported by the German Federal Ministry of Education and Research (BMBF), “Unternehmen Region”, instrument Wachstumskern-Potenzial (no. 03WKP45). More information: <http://www.digitale-kuratierung.de>.

## References

- Babakar, M., Moy, W.: The State of Automated Factchecking - How to make factchecking dramatically more effective with technology we have now. Full Fact (2016). [https://fullfact.org/media/uploads/full\\_fact-the\\_state\\_of\\_automated\\_factchecking\\_aug\\_2016.pdf](https://fullfact.org/media/uploads/full_fact-the_state_of_automated_factchecking_aug_2016.pdf)
- Barthel, M., Mitchell, A., Holcomb, J.: Many Americans Believe Fake News Is Sowing Confusion. Pew Research Center (2016). <http://www.journalism.org/2016/12/15/many-americans-believe-fake-news-is-sowing-confusion/>
- Bazzaz, D.: News you can use: infographic walks you through 10 questions to detect fake news. The Seattle Times (2016). <http://www.seattletimes.com/education-lab/infographic-walks-students-through-10-questions-to-help-them-spot-fake-news/>
- Belhajjame, K., Cheney, J., Corsark, D., Garijo, D., Soiland-Reyes, S., Zednik, S., Zhao, J.: PROV-O: The PROV Ontology. W3C Recommendation, World Wide Web Consortium (W3C), April 2013a. <https://www.w3.org/TR/2013/REC-prov-o-20130430/>
- Belhajjame, K., Deus, H., Garijo, D., Klyne, G., Missier, P., Soiland-Reyes, S., Zednik, S.: PROV Model Primer. W3C Working Group Note, World Wide Web Consortium (W3C), April 2013b. <https://www.w3.org/TR/2013/NOTE-prov-primer-20130430/>
- Bourgonje, P., Moreno-Schneider, J., Nehring, J., Rehm, G., Sasaki, F., Srivastava, A.: Towards a platform for curation technologies: enriching text collections with a semantic-web layer. In: Sack, H., Rizzo, G., Steinmetz, N., Mladenić, D., Auer, S., Lange, C. (eds.) ESWC 2016. LNCS, vol. 9989, pp. 65–68. Springer, Cham (2016a). [https://doi.org/10.1007/978-3-319-47602-5\\_14](https://doi.org/10.1007/978-3-319-47602-5_14)

- Bourgonje, P., Schneider, J.M., Rehm, G., Sasaki, F.: Processing document collections to automatically extract linked data: semantic storytelling technologies for smart curation workflows. In: Gangemi, A., Gardent, C. (eds.) Proceedings of the 2nd International Workshop on Natural Language Generation and the Semantic Web (WebNLG 2016), Edinburgh, UK, pp. 13–16. The Association for Computational Linguistics, September 2016b
- Bourgonje, P., Schneider, J.M., Rehm, G.: Automatic classification of abusive language and personal attacks in various forms of online communication. In: Rehm, G., Declerck, T. (eds.) GSCL 2017. LNAI, vol. 10713, pp. 180–191. Springer, Heidelberg (2017a)
- Bourgonje, P., Schneider, J.M., Rehm, G.: From clickbait to fake news detection: an approach based on detecting the stance of headlines to articles. In: Popescu, O., Strapparava, C. (eds.) Proceedings of the Second Workshop on Natural Language Processing meets Journalism - EMNLP 2017 Workshop (NLP MJ 2017), Copenhagen, Denmark, pp. 84–89, 7 September 2017b
- Chan, R.: Artificial Intelligence is Going to Destroy Fake News - But A.I. can also cause the volume of fake news to explode. Inverse Innovation (2017). <https://www.inverse.com/article/27723-artificial-intelligence-will-destroy-fake-news>
- Conroy, N.J., Rubin, V.L., Che, Y.: Automatic deception detection: methods for finding fake news. Proc. Assoc. Inf. Sci. Technol. **52**(1), 1–4 (2015)
- Crystal, D.: Language and the Internet. Cambridge University Press, Cambridge (2001)
- Dale, R.: NLP in a post-truth world. Nat. Lang. Eng. **23**(2), 319–324 (2017)
- Derczynski, L., Bontcheva, K., Liakata, M., Procter, R., Hoi, G.W.S., Zubiaga, A.: SemEval-2017 task 8: RumourEval: determining rumour veracity and support for rumours. In: Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), Vancouver, Canada, pp. 69–76. Association for Computational Linguistics, August 2017
- Ding, T., Bickel, W.K., Pan, S.: Social Media-based Substance Use Prediction, May 2017. <https://arxiv.org/abs/1705.05633>
- Filloux, F.: Quality for news is mostly about solving the reputation issue (2017). <https://mondaynote.com/quality-for-news-is-mostly-about-solving-the-reputation-issue-fdebd0dcc9e2>
- Gershgorn, D.: In the fight against fake news, artificial intelligence is waging a battle it cannot win. Quartz (2016). <https://qz.com/843110/can-artificial-intelligence-solve-facebooks-fake-news-problem/>
- Groth, P., Moreau, L.: PROV-overview: an overview of the PROV family of documents. W3C Working Group Note, World Wide Web Consortium (W3C), April 2013. <https://www.w3.org/TR/prov-overview/>
- Gupta, A., Lamba, H., Kumaraguru, P., Joshi, A.: Faking Sandy: characterizing and identifying fake images on Twitter during Hurricane Sandy. In: Proceedings of the 22nd International Conference on World Wide Web, Rio de Janeiro, pp. 729–736 (2013)
- Haase, M., Huber, M., Krumeich, A., Rehm, G.: Internetkommunikation und Sprachwandel. In: Weingarten, R. (ed.) Sprachwandel durch Computer, pp. 51–85. Westdeutscher Verlag, Opladen (1997)
- Holan, A.D.: 2016 lie of the year: fake news. PolitiFact (2016). <http://www.politifact.com/truth-o-meter/article/2016/dec/13/2016-lie-year-fake-news/>
- Horne, B.D., Adal, S.: This just in: fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news. In: Proceedings of the 2nd International Workshop on News and Public Opinion at ICWSM, March 2017

- Mantzarlis, A.: 6 tips to debunk fake news stories by yourself. Poynter (2015). <http://www.poynter.org/2015/6-tips-to-debunk-fake-news-stories-by-yourself/385625/>
- Mantzarlis, A.: There are now 114 fact-checking initiatives in 47 countries. Poynter (2017). <https://www.poynter.org/2017/there-are-now-114-fact-checking-initiatives-in-47-countries/450477/>
- Marchi, R.: With Facebook, blogs, and fake news, teens reject journalistic “Objectivity”. *J. Commun. Inq.* **36**(3), 246–262 (2012)
- Martinez-Alvarez, M.: How can machine learning and AI help solving the fake news problem? (2017). <https://miguelmalvarez.com/2017/03/23/how-can-machine-learning-and-ai-help-solving-the-fake-news-problem/>
- Marwick, A., Lewis, R.: Media Manipulation and Disinformation Online. Data and Society Research Institute, May 2017. <https://datasociety.net/output/media-manipulation-and-disinfo-online/>
- Metz, C.: The Bittersweet Sweepstakes to Build an AI that Destroys Fake News. *Wired* (2016). <https://www.wired.com/2016/12/bittersweet-sweepstakes-build-ai-destroys-fake-news/>
- Rehm, G.: Towards automatic web genre identification - a corpus-based approach in the domain of academia by example of the academic’s personal homepage. In: Proceedings of the 35th Hawaii International Conference on System Sciences (HICSS-35), Big Island, Hawaii. IEEE Computer Society, January 2002
- Rehm, G., Sasaki, F.: Digitale Kuratierungstechnologien - Verfahren für die effiziente Verarbeitung, Erstellung und Verteilung qualitativ hochwertiger Medieninhalte. In: Proceedings of the 2015 International Conference of the German Society for Computational Linguistics and Language Technology, GSCL 2015, pp. 138–139 (2015)
- Rehm, G., Sasaki, F., Burchardt, A.: Web Annotations - A Game Changer for Language Technologies? Presentation given at I Annotate 2016, Berlin, Germany, 19/20 May 2016. <http://www.slideshare.net/georgrehm/web-annotations-a-game-changer-for-language-technology>, <http://iannotate.org/2016/>
- Rehm, G., Schneider, J.M., Bourgonje, P., Srivastava, A., Nehring, J., Berger, A., König, L., Rächle, S., Gerth, J.: Event detection and semantic storytelling: generating a travelogue from a large collection of personal letters. In: Caselli, T., Miller, B., van Erp, M., Vossen, P., Palmer, M., Hovy, E., Mitamura, T. (eds.) Proceedings of the Events and Stories in the News Workshop, Vancouver, Canada. Association for Computational Linguistics. Co-located with ACL 2017, August 2017, in print
- Reid, E.M.: Electropolis: communication and community on internet relay chat. Honours thesis, University of Melbourne, Department of History (1991). <http://www.aluluei.com/electropolis.htm>
- Rogers, K., Bromwich, J.E.: The Hoaxes, Fake News and Misinformation We Saw on Election Day. *New York Times* (2016). <https://www.nytimes.com/2016/11/09/us/politics/debunk-fake-news-election-day.html>
- Rubin, V.L., Chen, Y., Conroy, N.J.: Deception detection for news: three types of fakes. *Proc. Assoc. Inf. Sci Technol.* **52**(1), 1–4 (2015). <https://doi.org/10.1002/pr2.2015.145052010083/full>
- Runkehl, J., Schlobinski, P., Siever, T.: Sprache und Kommunikation im Internet - Überblick und Analysen. Westdeutscher Verlag, Opladen, Wiesbaden (1998)
- Sanderson, R.: Web Annotation Protocol. W3C Recommendation, World Wide Web Consortium (W3C), February 2017. <https://www.w3.org/TR/2017/REC-annotation-protocol-20170223/>
- Sanderson, R., Ciccarese, P., Young, B.: Web Annotation Data Model. W3C Recommendation, World Wide Web Consortium (W3C), February 2017a. <https://www.w3.org/TR/2017/REC-annotation-model-20170223/>



- Sanderson, R., Ciccarese, P., Young, B.: Web Annotation Vocabulary. W3C Recommendation, World Wide Web Consortium (W3C), February 2017b. <https://www.w3.org/TR/2017/REC-annotation-vocab-20170223/>
- Schmidt, A., Wiegand, M.: A survey on hate speech detection using natural language processing. In: Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media, Valencia, Spain (2017)
- Srivastava, A., Rehm, G., Schneider, J.M.: DFKI-DKT at SemEval-2017 Task 8: rumour detection and classification using cascading heuristics. In: Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017), Vancouver, Canada, pp. 477–481. Association for Computational Linguistics, August 2017
- Walbrühl, D.: Das musst du wissen, um Fake News zu verstehen. Perspective Daily (2017). <https://perspective-daily.de/article/213/AhopoOEF>
- Wardle, C.: Fake news. It's complicated. First Draft News (2017). <https://firstdraftnews.com/fake-news-complicated/>
- Weedon, J., Nuland, W., Stamos, A.: Information Operations and Facebook, April 2017. Version 1.0 <https://fbnewsroomus.files.wordpress.com/2017/04/facebook-and-information-operations-v1.pdf>
- Zubiaga, A., Liakata, M., Procter, R., Hoi, G.W.S., Tolmie, P.: Analysing how people orient to and spread rumours in social media by looking at conversational threads. PLoS ONE 11(3), e0150989 (2016). <https://doi.org/10.1371/journal.pone.0150989>

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

