

12

Semantic Storytelling

From Experiments and Prototypes to a Technical Solution

Georg Rehm, Karolina Zaczynska, Peter Bourgonje,
Malte Ostendorff, Julián Moreno-Schneider, Maria Berger,
Jens Rauenbusch, André Schmidt, Mikka Wild, Joachim Böttger,
Joachim Quantz, Jan Thomsen, and Rolf Fricke

Abstract. In the past we experimented with variations of an approach we call semantic storytelling, in which we use multiple text analytics components including named entity recognition and event detection. This chapter summarises some of our previous work with an emphasis on the detection of movement action events and describes the long-term semantic storytelling vision as well as the setup and approach of our future work towards a robust technical solution, which is primarily driven by three industry use cases. Ultimately, we plan to contribute an implemented approach for semantic storytelling that makes use of various analytics services and that can be deployed in a flexible way in various industrial production environments.

12.1 Introduction: Technologies for Content Curation

With the ever increasing amount of digital content, users face the challenge of coping with enormous quantities of information. This is especially true for digital content curators; i.e., analysts or knowledge workers such as journalists (Moreno-Schneider et al., 2017b), television producers (Rehm et al., 2018), designers (Rehm et al., 2017a), librarians (Neudecker and Rehm, 2016), and academics (Rehm et al., 2019a), among others. These, and other, professional profiles have in common that they monitor and process *incoming content* with the goal of producing *new content*. This involves various processes – for example, to scan, translate, skim, contextualise, sort, summarise, evaluate, validate, cross-reference and to assess content – often under extreme time pressure. In several research and innovation projects (Rehm et al., 2015, 2019a) and most recently in QURATOR (Rehm et al., 2020b), we have been developing technical approaches with the goal of supporting knowledge

workers in their day-to-day jobs, curating large amounts of content with language and knowledge technologies more efficiently and more effectively. One of our focus areas is the identification and generation of storylines (see below). This includes both the (semi-)automatic creation of new content as well as helpful presentation and visualisation techniques that make use of the results of our semantic technologies. We call the approach semantic storytelling (Bourgonje et al., 2016b; Moreno-Schneider et al., 2016, 2017b; Rehm et al., 2017b, 2018, 2019b).

We focus on storytelling as a generic human technique to order a series of events in the world or, in a more abstract way, pieces of information, and find meaningful patterns in them (Bruner, 1991). By telling a story, we partially or fully relate events to a schematic structure – e.g., in terms of topic, locality or causal relationships – and construct explanatory models. Humans are able to dynamically adjust their narratives and tell their stories differently depending on who the listener is (Rishes et al., 2013), whereas this is still a challenging task for machines.

In recent years, storytelling has mostly been interpreted as a natural language generation (NLG) task (see, e.g., Fan et al., 2018, 2019), where the goal is actually to generate texts based on a headline or keywords, for example. We interpret the concept differently by concentrating on the *extraction* and *presentation* of stories and their individual parts, contained in incoming content streams; for example, document collections or social media feeds. The goal is to enable content curators to create new storylines based on the information extracted and presented by our technologies.

We see storylines as sets of building blocks that depending on their combination (temporal, geographical, causal etc.), can be assembled into a story in various ways. Our goal is the recognition of various atomic pieces of information (e.g., facts, propositions, entities, events) in multiple documents and the identification of semantic relations between these atomic pieces. Corresponding applications can be conceptualised, among others, as information systems (for the retrieval of existing content) or recommender systems (focusing upon the creation of new content).

The remainder of this chapter is structured as follows. First, Section 12.2 summarises our previous work on the topic, focusing on technical components and developments. Section 12.3 briefly presents our three current industry application scenarios from which we derive a set of requirements. These inform our most recent technical approach, which is especially geared towards robustness and flexibility as well as practical application in the three

use cases (Section 12.4). Section 12.5 presents related work. Section 12.6 concludes the chapter.

12.2 Semantic Storytelling: Selected Components

Semantic storytelling can be conceptualised as the automatic or semi-automatic generation of different storylines based on information extracted from extensive document collections or social media streams and then processed, classified, annotated and visualised, typically in an interactive way. In the following, we describe selected preliminary approaches and experimental components, as well as their interactions, that we developed over the years working on the topic.

12.2.1 Text and Document Analytics and Linked Data

A set of text analytics services is the technological foundation of our semantic storytelling architecture (see, e.g., Bourgonje et al., 2016a,b; Moreno-Schneider et al., 2016; Rehm et al., 2018). These belong to the following three larger groups:

1. Services that analyse complete documents (or document collections) to provide document-level metadata: language identification, document structure analysis, text genre detection, topic detection
2. Services that extract, annotate and enrich specific parts of the incoming content: named entity recognition (including rudimentary co-reference resolution), named entity linking, time expression analysis, topic detection, event detection
3. Services that transform parts of the content or whole documents: single document summarisation, multidocument summarisation, automated translation

The services are in different stages of maturity. They are orchestrated using a workflow manager and additional platform-related tooling (Moreno-Schneider et al., 2020). Eventually, they will be made available through the European Language Grid (Rehm et al., 2020a).

Named entity recognition and linking as well as time expression analysis are performed to identify named entities of various types and classes

(persons, locations, organisation etc.). Whenever possible, entities, topics etc. are anchored to external knowledge graphs (e.g., DBPedia, Wikidata, Geonames), which enables us to perform SPARQL queries to retrieve additional information for each item such as, e.g., Global Positioning System (GPS) coordinates for locations. The integration of the results of the time expression analysis allows reasoning over temporal expressions and anchoring entities and events to a timeline. We use topic detection to assign abstract topics to individual sentences, paragraphs, chapters and documents. Annotated topics constitute yet another layer of accessing and recombining the processed content. For the annotations themselves we use NLP Interchange Format (Hellmann et al., 2013), which allows the exploitation of the Semantic Web and Linked Data paradigm (including its vast set of formats, formalisms and tools) and Linked Open Data resources for storyline generation.

To be able to analyse a wide variety of incoming documents with the same setup, we distinguish between different classes or genres of documents; i.e., we experiment with different approaches for identifying document structures (Rehm et al., 2019b) and document genres (Rehm, 2007). An ontology to represent a heterogeneous set of document characteristics, essentially tying together all the different annotations mentioned above, is currently under development.

12.2.2 Detection of Movement Action Events

To describe one service and experiment in more detail, we implemented an event detection system based on Yang and Mitchell (2016) to pinpoint words or phrases in sentences that refer to events involving participants and locations affected by other events and spatiotemporal aspects. The module is trained on the ACE 2004 data (Doddington et al., 2004). This service can also perform event detection crosslingually by automatically translating non-English-language documents to English first and then detecting events in the translated documents.

The prototype described below enables putting together a story interactively based on semantically enriched content. As a first use case, we concentrated on the approximately 2,800 letters exchanged between German architect Erich Mendelsohn and his wife, Luise, between 1910 and 1953. The collection contains 2,796 letters, written between 1910 and 1953, with a total of 1,002,742 words (359 words per letter on average) on more than 11,000 sheets of paper. Most are in German (2,481); the rest is written in English (312) and French (3).

Table 12.1. *Automatically extracted movement action events (MAEs)*

Letter Text	Extracted MAEs
Another train stopped [...] this would be the train with which Eric had to leave Cleveland.	Eric, Cleveland, [], [], [], train
because I have to leave on the 13th for Chicago.	I [Erich], Croton on Hudson, NY, Chicago, 13th Dec. 1945, [], []
April 5th 48 Sweetheart – Here I am – just arrived in Palm Springs [...]	I [Erich], [], Palm Springs, [], 5th April 1948, []
Germaine wants me to come up with her to Tahoe – she will be there for 2 weeks from July 15th [...]	Germaine, [], Tahoe, [], 15th July, []
Thompsons are leaving for a week – [...] at the Beverly Hills on Thursday night!!	Thompsons, [], Beverly Hills, 8th July, [], []

In the letters, the Mendelsohns discuss their private and professional lives, their relationship, meetings with friends and business partners, and also their travels. Bienert and de Wit (2014) provide transcriptions of the letters together with scans, photos and metadata.

We want to transform this set of interconnected *letters* into a *travelogue* that provides an engaging story to the reader and that also enables additional modes of access; e.g., through map-based or timeline-based visualisations. In this experiment we explore to what extent it is possible to automate the production of a travelogue from a collection of letters. We focused on a specific class of events, movement action events (MAEs). A complete description of the experiment can be found elsewhere (Rehm et al., 2017b); here we only present a few examples of extracted MAEs to demonstrate the functionality (Table 12.1). In the letters, MAEs are typically mentioned whenever the author is undertaking or about to undertake a trip from A to B using a specific mode of transport. An MAE consists of the six-tuple $MAE = \langle P, L_O, L_D, t_d, t_a, m \rangle$, where P is a reference to the participant (E. or L. Mendelsohn), L_O and L_D are the origin and destination locations (named locations, GPS coordinates), t_d and t_a are the time of departure and arrival and m is the mode of transport. Each component is optional as long as the MAE contains a participant and a destination. We have been able to successfully retrace the various journeys of the Mendelsohns. Through additional modes of access to the content base (i.e., the letters themselves, photographs, sketches), by means

of maps or timelines, the authoring environment (Section 12.2.3) provides suggestions – i.e., potential story paths – to the content curator who is using the tool for putting together a new story on the Mendelsohns' lives and journeys.

12.2.3 Storytelling Prototypes and User Interfaces

Following the general approach described above, we implemented a number of experimental prototypes and corresponding interactive user interfaces (UIs; Bourgonje et al., 2016b; Rehm et al., 2017b; Moreno-Schneider et al., 2016; Rehm et al., 2018). On top of the semantic analysis of document collections, we map the extracted information, whenever possible, to Linked Open Data (LOD), and visualise the result (Moreno-Schneider et al., 2017a; Rehm et al., 2017a). The UIs support content curators in the process of identifying interesting or surprising relationships between different concepts or entities mentioned in the processed documents. By providing feedback to the output of certain semantic services, content curators have some amount of control on the workflow. They are able to upload existing language resources to adapt individual services. For example, the Named Entity Recognition (NER) service allows users to supply dictionaries for entity linking and the event detection service allows users to supply lists of entities for the identification of agents for events.

The storytelling UIs involve the dynamic and interactive recomposition and visualisation of extracted information based on the information extracted from the text analytics services. This especially involves arranging content elements (documents, paragraphs, sentences, claims or events) on a dynamic timeline. The summarisation services are used to compress larger pieces of content into bites that can be easily digested, moved around on the screen and maybe expanded back into their original versions. We currently experiment with the algorithmic construction of storylines based on the recomposition of previously extracted information (for more details, see Section 12.4).

To complement the experiment in which we extracted movement action events (Section 12.2.2), an authoring environment was developed, and several screens of the UI are shown in Figure 12.1. It was a conscious design decision to move beyond the typical notion of a 'web page' that is broken up into different 'modules' using templates. The focus of this prototype is the development of engaging stories told through informative content. With this tool the content curator can interactively put together a story based on the semantically enriched content. In this example case we work with the letters exchanged between Erich and Luise Mendelsohn (see Section 12.2.2).

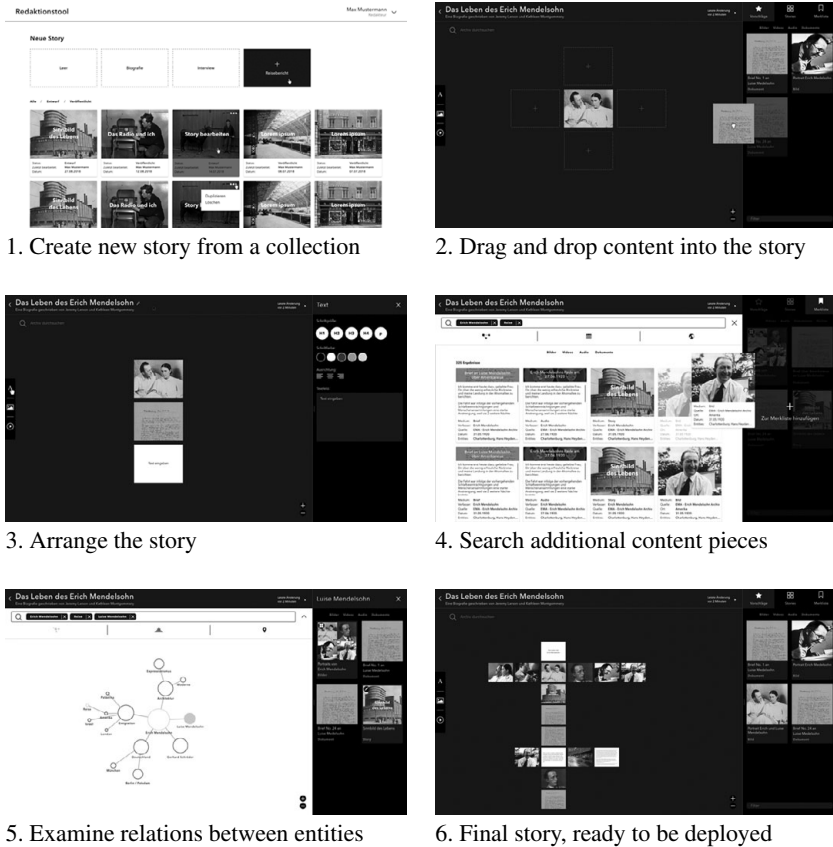


Figure 12.1 Semantic storytelling authoring environment.

12.3 Semantic Storytelling in Industry Use Cases

The emerging semantic storytelling technologies described in this chapter, especially in Section 12.4, are primarily developed for three of the industry partners in the QURATOR project (Rehm et al., 2020b). The goal is to integrate them as base technologies into the three use cases (UC1, UC2, UC3) and prototype applications. In the following subsections we describe these use cases, from which we derive a set of basic requirements that inform the final design of our technical approach.

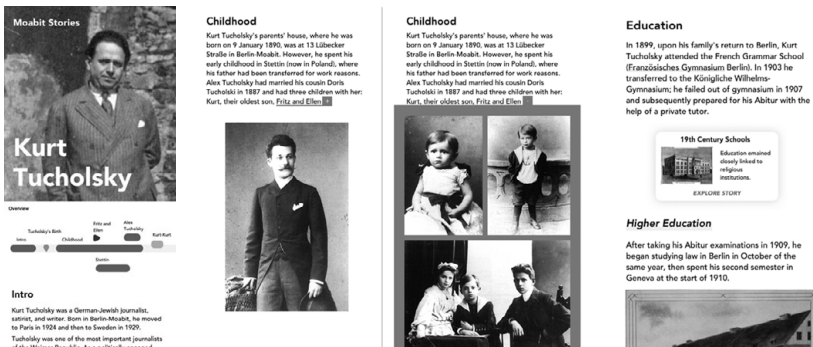


Figure 12.2 Example screens of the ‘Explore the Neighbourhood!’ concept.

12.3.1 Use Case: Explore the Neighbourhood!

The first use case, conceptualised by 3pc GmbH (<http://3pc.de>), enables urban explorers to discover a neighbourhood (UC1). The idea is to collect content from online sources, cultural heritage archives or publishing houses to identify relevant content that matches the current location. This enables the location-aware application to ‘tell stories’ about the neighbourhood. In the prototype we use Moabit, a multicultural Berlin district with a rich history and lively present. Our vision for the app is to present stories that have been generated (semi-)automatically according to their location, time, relevance for a certain user, location and budget.

The app presents curated interactive stories to a user who is located in Moabit or generates new ones. Figure 12.2 shows an example of the option to explore the life of Kurt Tucholsky, a German-Jewish journalist and writer, born in Moabit in 1890. The user can follow the linear story through Tucholsky’s life, consuming texts, (historic) images, video and audio files and interactive elements. Stories are conceptualised as hypertexts, where users can dive deeper into certain aspects, expanding, shortening or bypassing topics. The app also reacts to data, such as weather, traffic, and, crucially, the current location. It has access to further information, such as opening times of points of interest. Augmenting these pieces of information through digitised archives enables stories with interactive paths and junctions that are better suited to the users’ personal interests. Semantic storytelling is needed to support editors in their curation process and, eventually, to generate stories automatically.

12.3.2 Use Case: Smart Newsboard

The Smart Newsboard is a vision by Condat AG (<http://condat.de>) that uses semantic storytelling technologies for the production of new content based on articles on a specific topic (UC2). This involves various curation services, such as, e.g., (1) finding original sources (through targeted searches in RSS feeds or social media); (2) categorising these sources into clusters through topic detection or text classification; (3) applying text analytics and recognising and linking named entities to resources such as Wikidata. This allows the creation of knowledge graphs to identify connections between people and events, to enable journalists to take a deep dive. A crucial feature is the detection of temporal expressions, assigning them to facts or claims, whether they are absolute (e.g., *2019, October 3rd*) or relative (e.g., *last week, next Monday, after that*), that can be used to generate timelines. Other important services include fact checking to help journalists evaluate the trustworthiness of content and multidocument summarisation (Aksenov et al., 2020). Putting these services in practice, the outline of a news story can be generated together with related content. The main focus of storytelling in the Smart Newsboard is to generate a story outline in a linear and neutral way.

12.3.3 Use Case: Wikipedia Trails

This use case, devised by ART+COM AG (<http://artcom.de>), produces a visual map based on Wikipedia usage, with the aim of supporting knowledge workers and analysts in the process of acquiring, recording, and passing on knowledge (UC3; Figure 12.3). Inspired by Vannevar Bush's seminal essay

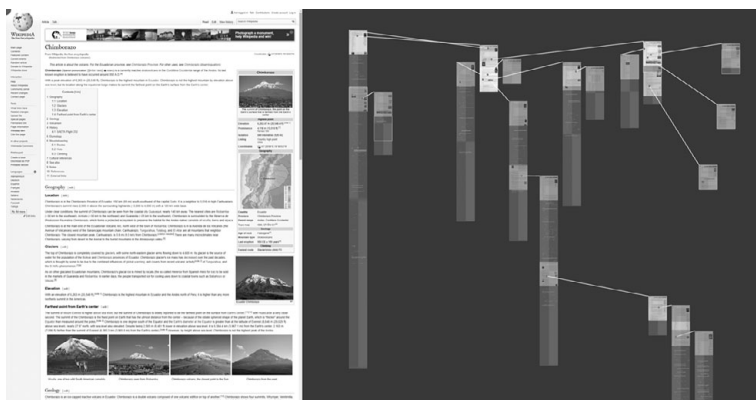


Figure 12.3 Example screen of the 'Wikipedia Trails' concept.

‘As We May Think’ (Bush, 1945), the tool records timestamped snapshots of Wikipedia pages while a user navigates through them. Positions in pages and clicked links are also recorded. The data are continually pieced together in a visualisation, yielding long, thumbnail-like representations of pages. These thumbnails are then placed in relation to each other using force-directed graph layout techniques. The approach allows for the gradual buildup of what we imagine Bush termed a ‘trail’: an integrated recording of how users navigate a knowledge space. The aim is a novel method that sits somewhere between existing perspectives for the preservation of knowledge acquisition that are either too broad (browser history) or too piecemeal (files, bookmarks, tabs) by visualising a map of the process itself.

12.3.4 Integrating Semantic Storytelling

From the three use cases, we can derive requirements for the methodology of the prototype development. Our goal is to design the approach in such a way that we can implement, based on the various experimental developments presented in Section 12.2, *one* robust and flexible technology solution that addresses the requirements of *all three* industry project partners at the same time.

UC1 needs storytelling support for the (semi-)automatic production of interactive stories. Either a human editor or an automated system is putting together stories that can be experienced by users of the mobile app. UC2 is similar to UC1. In both use cases, semantic storytelling functionalities are needed to suggest content relevant for the topic on which a journalist is currently writing a story. In contrast to UC1, in UC2 this needs to happen dynamically and in real time. The main purpose of UC3 is the dynamic visualisation of the browsing history of Wikipedia pages to support knowledge workers in desk research and knowledge acquisition tasks; i.e., storytelling technologies are meant to provide further insights into the semantic relationships between Wikipedia pages beyond the mere fact that hyperlinks exist between them.

Regarding the input, all three use cases need support for the processing of complete documents but also for addressing (including linking to) only parts of documents – i.e., individual paragraphs or sentences – to expose only key sentences to the user. The technology needs to be able to handle different text types, from historical cultural heritage documents and information made available by a municipality or city administrations (UC1) to RSS feeds and news articles (UC2) to web texts, e.g., Wikipedia pages (UC3). All three include multimedia elements, the processing of which is also needed. Furthermore, UC1 and UC2 have a demand for services such as summarisation and machine translation.

12.4 Towards a Flexible and Robust Technology Solution for Semantic Storytelling

We performed an analysis of the three use cases with the goal of identifying technical building blocks that can be shared by UC1, UC2 and UC3. The key insight of this analysis is that three simple yet technically ambitious processing steps, operating on two text segments, are needed to address the demands in the three use cases, in addition to the various tools and services mentioned previously (Section 12.2).

We can illustrate the functionality of these three technical building blocks with an example: a journalist is working on a story on a topic *T*. Our goal is to identify and to suggest new content that can be included in the emerging story. First, we need to identify whether, e.g., a certain tweet or text from an incoming RSS feed is *relevant* for *T*. Second, we need to determine the *importance* of the incoming tweet or article for the emerging story. Third, we identify the *semantic relation* between the incoming and the emerging article, which could be, among others, *background*, *cause*, *contrast* and *example*. The result of the third step can be also used in terms of visualisation; e.g., putting background information in the margins or recommending arguments in favour of or against a certain point of view; i.e., the approach should be able to extract and classify semantic relations between text segments (see Figure 12.4 and Rehm et al., 2019c, for more details).

Below, we describe the three steps (identification of a new or incoming text segment's *relevance* and *importance* as well as the *semantic or discourse relation* that holds between the new, incoming segment and topic *T*) that are at the core of our approach in more detail. Note that steps 1 and 2 as well as steps 2 and 3 overlap in terms of their corresponding scope. For example, when dealing with self-contained document collections in UC1, step 1 can be omitted because *relevance* may be an inherent property of a collection that includes content about a certain neighbourhood. Similarly, the type of relation that holds between two segments can also bear information on their relative importance, rendering the separation between steps 2 and 3 less strict.

Step 1: Determine the Relevance of a Segment for a Topic The approach starts with a topic *T*, instantiated through a text segment (e.g., a complete document, headline, summary, or a named entity), specified by a journalist who is writing a story on *T*. In order to identify content pieces relevant for *T*, we process – maybe even continuously – incoming news feeds, self-contained document collections, systematically compiled corpora or knowledge bases. For each piece of content processed, we need to decide whether its topic is relevant for *T*.

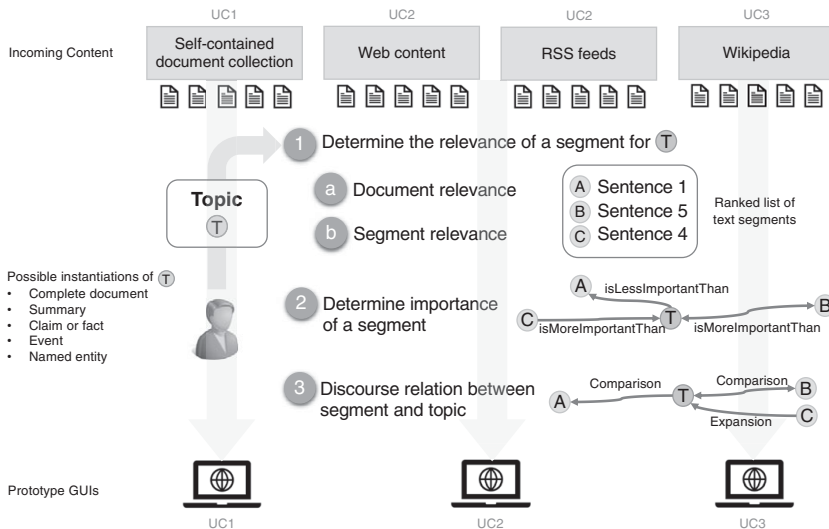


Figure 12.4 Architecture of the approach.

Relevance can be computed in various ways. We can employ text similarity measures, approaches from information retrieval or the overlap in terms of named entities contained in the segments. Crucially, the (accuracy of the) approach is dependent on the length of the segments. We currently concentrate on topic modeling (latent Dirichlet allocation [LDA], latent semantic analysis [LSA]). Without explicitly modeling topics, we can also perform pair-wise comparisons of document similarity (see, e.g., Salton and McGill, 1983; Pagliardini et al., 2018; Selivanov and Wang, 2016).

Step 2: Determine the Importance of a Segment If a document d is related to T , we can determine the importance of d (and its individual segments) with regard to T . Though this task is related to step 1, off-the-shelf approaches for determining the importance of a text segment with regard to a topic do not exist. Various cues can potentially be exploited, though. For example, with regard to UC2, an incoming news piece on T , published seconds ago, that includes the cue word ‘BREAKING’ in its title can surely be considered important. Determining the topical importance can also be framed as a question answering task, where T is the question. Transformer-based language models have achieved state-of-the-art results for question answering (Devlin et al., 2019), highlighting the relevance of these model architectures for storytelling.

We can also borrow from rhetorical structure theory (RST, Mann and Thompson, 1988). The construction of an RST tree involves various decisions

with regard to the status of text segments including their discourse relation to other segments and also regarding their role as a *nucleus* (N, the important part of a relation) or as a *satellite* (S, the contributing part) in a specific discourse relation. Two segments are assigned either an S–N, N–S or N–N structure. This subtask can be done in isolation (Hernault et al., 2010; Soricut and Marcu, 2003) or combined with the relation classification task (Joty et al., 2015). Performed iteratively, this pair-wise classification can result in a set of most important segments regarding *T*. Due to low amounts of training data (Carlson et al., 2002), resulting in less robust classifiers, we follow Sporleder and Lapata (2005) in isolating the nuclearity assignment task.

Step 3: Determine the Discourse Relation between Two Segments The modeling of coherence relations in textual content is at the core of discourse processing and corresponding parsing frameworks. These analyse a text for, typically, intratextual but intersentential relations. Some of the frameworks also indicate, given two segments, which is the more important one. We borrow from discourse parsing, but there are several added challenges. Crucially, our system needs to be able to robustly process and compare short segments (typically, sentences or paragraphs) extracted from *different* texts for which we have ample evidence from step 1 that the two texts are relevant to each other. After having established the relevance and relative importance, we proceed with determining the discourse or semantic relation that exists between these individual segments from *different* texts. Our initial experiments are based on the Penn Discourse Tree Bank (PDTB). We adopt Penn Discourse Tree Bank's (PDTB) four top-level senses and achieve promising results (Rehm et al., 2019c). Once established between the two segments, the key idea is to make use of the discourse relation in visualisations and UIs recommendations, text generation approaches etc.

12.5 Related Work

In the following, we provide a brief overview of related work conducted in several areas including narratology, discourse theory, and applied work in computational linguistics and language technology.

Several approaches grounded in narratology address storytelling as a way of automatising the detection of instances of story grammars (Rumelhart, 1975), especially events, in texts. Vossen et al. (2021) present a data set for

the detection of temporal and causal relations and use a plot structure (Bal, 1997) to order events found in narratives or, more generally, text documents, chronologically and logically. According to Bal (1997), narratives follow a plot structure that consists of ordered events, told by an agent or author and caused or experienced by actors. Yarlott and Finlayson (2016) use Propp's (1968) morphology of Russian hero tales for story detection and generation systems. In his book, first published in 1928, Propp analyzes the basic structural elements of Russian folk tales, which always occur in a fixed, consecutive order.

Yan et al. (2019) describe a system that learns 'functional story schemas' as sets of functional structures (e.g., character introduction, conflict setup, etc.) in social media narratives. They extract patterns of functional structures. Afterwards their formation in a story is analyzed across all stories to find schematic structures. In contrast, Gordon et al. (2011) use stories from blog articles to perform automated causal reasoning. Their statistical and information retrieval-focused experiments show that a simple co-occurrence measure among the words of an antecedent (cause, reason) and a consequent (result) in a corpus of personal stories can help to derive causal information.

Cucchiarelli et al. (2018) examine news recommendations and present a technique that, for a given event, suggests news articles not covered previously by a journalist's research work. They measure an event in relation to its media coverage and compare it to the potential reader's communicative patterns, according to Twitter and Wikipedia analyses. The method recommends topics of interest to the reader that are only poorly covered in published news but emerge as topics of interest in Twitter and Wikipedia. In contrast, Bois et al. (2017) recommend articles based on simple lexical similarity. They link news articles in the form of a graph and label links to inform users on the nature of the relation between two news pieces.

Ribeiro et al. (2017) cluster news articles based on identified event instances and word alignment. They attempt to form clusters of online articles that deal with a certain event type. Nie et al. (2019) use dependency parsing and discourse relations to determine sentence relations by learning corresponding vector representations. Yarlott et al. (2018) apply the theory of hierarchical discourse by Dijk (1988) to examine how paragraphs behave when used as discourse structure units in news articles. They describe the relations of paragraphs with regard to the events in an article. Learning the discourse structure of sections of news articles helps the authors to understand the importance and temporal order of story items.

12.6 Conclusions and Future Work

After laying the groundwork for the full implementation of semantic storytelling through various experiments in a number of areas, we are now approaching the last phase, in which we attempt to combine the different preliminary pieces described in Section 12.2, including various text analytics and semantic enrichment services that operate on documents, document collections or smaller segments of documents (Rehm et al., 2019b), with the emerging set of technologies described in Section 12.4. To this end, we identified the key common building blocks of three industry use cases for semantic storytelling technologies. Though step 1 can be implemented using one of several known approaches, steps 2 and 3 are much more challenging (Figure 12.4). Essentially, our approach is grounded in the assumption that different texts that deal with the *same* topic but that are from *different* authors and *different* sources can be interconnected in a meaningful way through discourse relations, which we attempt to extract automatically in order to expose the identified relations to the respective downstream application. It is exactly these semantic relations or discourse relations and their directionality that we want to expose to the respective user in the corresponding prototype. With regard to the authoring environment, such relations hold between two media assets (as seen in the screenshots shown in Figure 12.1), and we want to support content curators making use of these relations by exposing them explicitly and exploiting them in the construction of storylines in a semiautomatic or fully automatic way. Though our initial experiments are promising (Rehm et al., 2019c), they also show that additional research is needed before we can integrate these technologies into the respective prototypes. Data sets annotated for rhetorical structure, discourse structure or, closely related, event structure are still rather limited both in availability and in size.

Our future work will thus focus on expanding our setup significantly, especially with regard to the analysis and classification of discourse relations and more sophisticated processing of connectives. We will also integrate a more flexible approach with regard to the processing of single documents by concentrating on larger parts of a document including longer summaries and paraphrased variants to increase coverage and recall. Taking into account explicit ontological knowledge to identify semantic relations between texts will also be an important next step towards the completion of the envisaged semantic storytelling architecture (Rehm et al., 2019b).

Acknowledgments This work is funded by the German Federal Ministry of Education and Research (BMBF) through the project QURATOR (<http://qurator.ai>, Unternehmen Region, Wachstumskern, grant no. 03WKDA1A).

References

- Aksenov, Dmitrii, Moreno-Schneider, Julián, Bourgonje, Peter, Schwarzenberg, Robert, Hennig, Leonhard, and Rehm, Georg. 2020. Abstractive Text Summarization Based on Language Model Conditioning and Locality Modeling. Pages 6682–6691 of: Calzolari, Nicoletta, Bächtel, Frédéric, Blache, Philippe, et al. (eds.), *Proceedings of the 12th Language Resources and Evaluation Conference (LREC 2020)*. Marseille, France: European Language Resources Association.
- Bal, Mieke. 1997. *Narratology: Introduction to the Theory of Narrative*. University of Toronto Press.
- Bienert, Andreas, and de Wit, Wim (eds.). 2014. *EMA – Erich Mendelsohn Archiv. Der Briefwechsel von Erich und Luise Mendelsohn 1910–1953*. Berlin: Staatliche Museen zu Berlin.
- Bois, Rémi, Gravier, Guillaume, Jamet, Eric, Robert, Maxime, Morin, Emmanuel, Sébillot, Pascale, and Robert, Maxime. 2017. Language-Based Construction of Explorable News Graphs for Journalists. Pages 31–36 of: Popescu, Octavian, and Strapparava, Carlo (eds.), *Proceedings of the 2017 EMNLP Workshop on Natural Language Processing Meets Journalism*. Copenhagen: Association for Computational Linguistics.
- Bourgonje, Peter, Moreno-Schneider, Julián, Rehm, Georg, and Sasaki, Felix. 2016a. Processing Document Collections to Automatically Extract Linked Data: Semantic Storytelling Technologies for Smart Curation Workflows. Pages 13–16 of: Gangemi, Aldo, and Gardent, Claire (eds.), *Proceedings of the 2nd International Workshop on Natural Language Generation and the Semantic Web (WebNLG 2016)*. Edinburgh: The Association for Computational Linguistics.
- Bourgonje, Peter, Moreno-Schneider, Julián, Nehring, Jan, Rehm, Georg, Sasaki, Felix, and Srivastava, Ankit. 2016b. Towards a Platform for Curation Technologies: Enriching Text Collections with a Semantic-Web Layer. Pages 65–68 of: Sack, Harald, Rizzo, Giuseppe, Steinmetz, Nadine, Mladenec, Dunja, Auer, Sören, and Lange, Christoph (eds.), *The Semantic Web*. Lecture Notes in Computer Science, no. 9989. Springer.
- Bruner, Jerome. 1991. The Narrative Construction of Reality. *Critical Inquiry*, **18**(1), 1–21.
- Bush, Vannevar. 1945. As we may think. *The Atlantic*, **176**(1), 101–108.
- Carlson, Lynn, Marcu, Daniel, and Okurowski, Mary Ellen. 2002. *RST Discourse Treebank*. Linguistic Data Consortium, Online. <https://catalog.ldc.upenn.edu/LDC2002T07>.
- Cucchiarelli, Alessandro, Morbidoni, Christian, Stilo, Giovanni, and Velardi, Paola. 2018. What to Write and Why: A Recommender for News Media. Pages 1321–1330 of: Haddad, Hisham M., Wainwright, Roger L., and Chbeir, Richard (eds.), *Proceedings of the 33rd Annual ACM Symposium on Applied Computing*. ACM.
- Dijk, Teun van. 1988. *News as Discourse*. Communication Series. Lawrence Erlbaum Associates.
- Doddington, George, Mitchell, Alexis, Przybocki, Mark, Ramshaw, Lance, Strassel, Stephanie, and Weischedel, Ralph. 2004. The Automatic Content Extraction

- (ACE) Program – Tasks, Data, and Evaluation. In: *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*. Lisbon: European Language Resources Association.
- Fan, Angela, Lewis, Mike, and Dauphin, Yann. 2018. Hierarchical Neural Story Generation. Pages 889–898 of: Gurevych, Iryna, and Miyao, Yusuke (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, Volume 1: Long Papers*. Association for Computational Linguistics.
- Fan, Angela, Lewis, Mike, and Dauphin, Yann. 2019. *Strategies for Structuring Story Generation*. arXiv preprint arXiv:1902.01109.
- Gordon, Andrew S., Bejan, Cosmin A., and Sagae, Kenji. 2011. Commonsense Causal Reasoning Using Millions of Personal Stories. In: Burgard, Wolfram, and Roth, Dan (eds.), *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*. AAAI Press.
- Hellmann, Sebastian, Lehmann, Jens, Auer, Sören, and Brümmer, Martin. 2013. Integrating NLP Using Linked Data. Pages 98–113 of: Alani, Harith, Kagal, Lalana, Fokoue, Achille, Groth, et al. (eds), *Proceedings of the 12th International Semantic Web Conference*. 21–25 October. Springer.
- Hernault, Hugo, Prendinger, Helmut, duVerle, David A., and Ishizuka, Mitsuru. 2010. HILDA: A Discourse Parser Using Support Vector Machine Classification. *Dialog & Discourse*, **1**(3), 1–33.
- Joty, Shafiq, Carenini, Giuseppe, and Ng, Raymond T. 2015. CODRA: A Novel Discriminative Framework for Rhetorical Analysis. *Computational Linguistics*, **41**(3), 385–435.
- Mann, William C., and Thompson, Sandra A. 1988. Rhetorical Structure Theory: Towards a Functional Theory of Text Organization. *Text-interdisciplinary Journal for the Study of Discourse*, **8**(3), 243–281.
- Mikolov, Tomas, Chen, Kai, Corrado, Greg, and Dean, Jeffrey. 2013. Efficient Estimation of Word Representations in Vector Space. In: Bengio, Yoshua, and LeCun, Yann (eds.), *1st International Conference on Learning Representations*.
- Moreno-Schneider, Julián, Bourgonje, Peter, Nehring, Jan, et al. 2016. Towards Semantic Story Telling with Digital Curation Technologies. In: Birnbaum, Larry, Popescu, Octavian, and Strapparava, Carlo (eds.), *Proceedings of Natural Language Processing Meets Journalism – IJCAI-16 Workshop (NLP MJ 2016)*.
- Moreno-Schneider, Julián, Bourgonje, Peter, Kintzel, Florian, and Rehm, Georg. 2020. A Workflow Manager for Complex NLP and Content Curation Pipelines. Pages 73–80 of: Rehm, Georg, Bontcheva, Kalina, Choukri, Khalid, Hajic, Jan, Piperidis, Stelios, and Vasiljevs, Andrejs (eds.), *Proceedings of the 1st International Workshop on Language Technology Platforms (IWLTP 2020, co-located with LREC 2020)*. 16 May.
- Moreno-Schneider, Julián, Bourgonje, Peter, and Rehm, Georg. 2017a. Towards User Interfaces for Semantic Storytelling. Pages 403–421 of: Yamamoto, Sakae (ed.), *Human Interface and the Management of Information: Information, Knowledge and Interaction Design, 19th International Conference, HCI International 2017 (Vancouver, Canada)*. *Lecture Notes in Computer Science (LNCS)*, no. 10274, Part II. Cham, Switzerland: Springer.
- Moreno-Schneider, Julián, Srivastava, Ankit, Bourgonje, Peter, Wabnitz, David, and Rehm, Georg. 2017b. Semantic Storytelling, Cross-Lingual Event Detection and

- Other Semantic Services for a Newsroom Content Curation Dashboard. Pages 68–73 of: *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing Meets Journalism*.
- Neudecker, Clemens, and Rehm, Georg. 2016. Digitale Kuratierungstechnologien für Bibliotheken. *Zeitschrift für Bibliothekskultur* 027.7, 4(2).
- Nie, Allen, Bennett, Erin, and Goodman, Noah. 2019. DisSent: Learning Sentence Representations from Explicit Discourse Relations. Pages 4497–4510 of: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Pagliardini, Matteo, Gupta, Prakhar, and Jaggi, Martin. 2018. Unsupervised Learning of Sentence Embeddings Using Compositional n-Gram Features. Pages 528–540 of: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Propp, Vladimir Y. 1968. *Morphology of the Folktale*. University of Texas Press (Original publication date 1928).
- Rehm, Georg. 2007. *Hypertextsorten: Definition – Struktur – Klassifikation*. PhD Thesis, Justus-Liebig-Universität Giessen.
- Rehm, Georg, Berger, Maria, Elsholz, Ela, et al. 2020a. European Language Grid: An Overview. Pages 3359–3373 of: *Proceedings of the 12th Language Resources and Evaluation Conference (LREC 2020)*. Marseille, France: European Language Resources Association.
- Rehm, Georg, Bourgonje, Peter, Hegele, Stefanie, et al. 2020b. QURATOR: Innovative Technologies for Content and Data Curation. In: Paschke, Adrian, Neudecker, Clemens, Rehm, Georg, Qundus, Jamal Al, and Pintscher, Lydia (eds.), *Proceedings of QURATOR 2020 – The Conference for Intelligent Content Solutions*. January 2021.
- Rehm, Georg, He, Jing, Moreno-Schneider, Julian, Nehring, Jan, and Quantz, Joachim. 2017a. Designing User Interfaces for Curation Technologies. Pages 388–406 of: Yamamoto, Sakae (ed.), *Human Interface and the Management of Information: Information, Knowledge and Interaction Design, 19th International Conference, HCI International 2017 (Vancouver, Canada). Lecture Notes in Computer Science (LNCS)*, no. 10273, Part I. Cham, Switzerland: Springer.
- Rehm, Georg, Lee, Martin, Schneider, Julián Moreno, and Bourgonje, Peter. 2019a. Curation Technologies for a Cultural Heritage Archive: Analysing and Transforming a Heterogeneous Data Set into an Interactive Curation Workbench. Pages 117–122 of: Antonacopoulos, Apostolos, and Büchler, Marco (eds.), *DATeCH 2019: Proceedings of the 3rd International Conference on Digital Access to Textual Cultural Heritage*. 8–10 May.
- Rehm, Georg, Moreno Schneider, Julián, Bourgonje, Peter, et al. 2017b. Event Detection and Semantic Storytelling: Generating a Travelogue from a large Collection of Personal Letters. Pages 42–51 of: Caselli, Tommaso, Miller, Ben, van Erp, Marieke, et al. (eds.), *Proceedings of the Events and Stories in the News Workshop*.
- Rehm, Georg, Moreno Schneider, Julián, Bourgonje, Peter, et al. 2018. Different Types of Automated and Semi-Automated Semantic Storytelling: Curation Technologies for Different Sectors. Pages 232–247 of: Rehm, Georg, and Declerck, Thierry

- (eds.), *Language Technologies for the Challenges of the Digital Age: 27th International Conference, GSCL 2017, Berlin, Germany, 13–14 September, 2017. Lecture Notes in Artificial Intelligence (LNAI)*, no. 10713. Springer.
- Rehm, Georg, and Sasaki, Felix. 2015. Digitale Kuratierungstechnologien – Verfahren für die eziante Verarbeitung, Erstellung und Verteilung qualitativ hochwertiger Medieninhalte. Pages 138–139 of: *Proceedings der Frühjahrstagung der Gesellschaft für Sprachtechnologie und Computerlinguistik (GSCL 2015)*. 30 September–2 October.
- Rehm, Georg, Zaczynska, Karolina, and Schneider, Julián Moreno. 2019b. Semantic Storytelling: Towards Identifying Storylines in Large Amounts of Text Content. Pages 63–70 of: Jorge, Alpio, Campos, Ricardo, Jatowt, Adam, and Bhatia, Sumit (eds.), *Proceedings of Text2Story – Second Workshop on Narrative Extraction From Texts co-located with 41th European Conference on Information Retrieval (ECIR 2019)*. 14 April.
- Rehm, Georg, Zaczynska, Karolina, Schneider, Julián Moreno, et al. 2020c. Towards Discourse Parsing-Inspired Semantic Storytelling. In: Paschke, Adrian, Neudecker, Clemens, Rehm, Georg, Qundus, Jamal Al, and Pintscher, Lydia (eds.), *Proceedings of Semantic Storytelling 269 QURATOR 2020 – The Conference for Intelligent Content Solutions*. 20–21 January.
- Ribeiro, Swen, Ferret, Olivier, and Tannier, Xavier. 2017. Unsupervised Event Clustering and Aggregation from Newswire and Web Articles. Pages 62–67 of: Popescu, Octavian, and Strapparava, Carlo (eds.), *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing Meets Journalism*. Copenhagen: Association for Computational Linguistics.
- Rishes, Elena, Lukin, Stephanie M., Elson, David K., and Walker, Marilyn A. 2013. Generating Different Story Tellings from Semantic Representations of Narrative. Pages 192–204 of: Koenitz, Hartmut, Sezen, Tonguc Ibrahim, Ferri, Gabriele, et al. (eds.), *Interactive Storytelling: 6th International Conference, ICIDS 2013*. Springer International Publishing.
- Rumelhart, David E. 1975. Notes on a Schema for Stories. Pages 211–236 of: *Representation and Understanding*. Elsevier.
- Salton, G., and McGill, M.J. 1983. *Introduction to Modern Information Retrieval*. Computer Series. New York: McGraw-Hill.
- Selivanov, Dmitriy, and Wang, Qing. 2016. text2vec: Modern Text Mining Framework for R. <https://cran.r-project.org/web/packages/text2vec/index.html>.
- Soricut, Radu, and Marcu, Daniel. 2003. Sentence Level Discourse Parsing Using Syntactic and Lexical Information. Pages 228–235 of: Hearst, Marti A., and Ostendorf, Mari (eds.), *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Sporleder, Caroline, and Lapata, Mirella. 2005. Discourse Chunking and Its Application to Sentence Compression. Pages 257–264 of: *HLT/EMNLP 2005, Human Language Technology Conference and Conference on Empirical Methods in*

- Natural Language Processing, Proceedings of the Conference, 6–8 October 2005, Vancouver, British Columbia, Canada.*
- Vossen, Piek, Caselli, Tommaso, and Segers, Roxane. 2021. A Narratology-Based Framework for Storyline Extraction. Pages 130–147 of: Caselli, Tommaso, Palmer, Martha, Hovy, Eduard, and Vossen, Piek (eds.), *Computational Analysis of Storylines: Making Sense of Events*. Cambridge University Press.
- Yan, Xinru, Naik, Aakanksha, Jo, Yohan, and Rose, Carolyn. 2019. Using Functional Schemas to Understand Social Media Narratives. Pages 22–33 of: *Proceedings of the Second Workshop on Storytelling*.
- Yang, Bishan, and Mitchell, Tom. 2016. Joint Extraction of Events and Entities within a Document Context. Pages 289–299 of: Knight, Kevin, Nenkova, Ani, and Rambow, Owen (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the ACL: Human Language Technologies*. Association for Computational Linguistics.
- Yarlott, W. Victor, Cornelio, Cristina, Gao, Tian, and Finlayson, Mark. 2018. Identifying the Discourse Function of News Article Paragraphs. Pages 25–33 of: Caselli, Tommaso, Miller, Ben, van Erp, Marieke, et al. (eds.), *Proceedings of the Workshop Events and Stories in the News 2018*. Santa Fe, NM: Association for Computational Linguistics.
- Yarlott, W. Victor H., and Finlayson, Mark A. 2016. ProppML: A Complete Annotation Scheme for Proppian Morphologies. Pages 8:1–8:19 of: Miller, Ben, Lieto, Antonio, Ronfard, Rémi, Ware, Stephen G., and Finlayson, Mark A. (eds.), *7th Workshop on Computational Models of Narrative (CMN 2016)*. OpenAccess Series in Informatics (OASISs), Vol. 53. Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum für Informatik.