**REPORT**RAPPORT

*INS*

Information Systems

*INformation Systems*

Composing discourse based on genre semantics

K.I. Falkovych, F.-M. Nack

CWI is the National Research Institute for Mathematics and Computer Science. It is sponsored by the Netherlands Organization for Scientific Research (NWO).
CWI is a founding member of ERCIM, the European Research Consortium for Informatics and Mathematics.

CWI's research has a theme-oriented structure and is grouped into four clusters. Listed below are the names of the clusters and in parentheses their acronyms.

Probability, Networks and Algorithms (PNA)

Software Engineering (SEN)

Modelling, Analysis and Simulation (MAS)

Information Systems (INS)

# Composing discourse based on genre semantics

ABSTRACT
The availability of semantically enriched data in repositories of larger content providers offers a means for new ways of multimedia presentation authoring. Existing (semi-)automatic content composition environments explore limited numbers of genres and concentrate on template-based approaches for composing discourse structures for these genres. We analyze a number of genres based on identifiable genre characteristics. We show that existing template-based approaches to content composition support the essay and biography genres that belong to one characteristic genre category but fail in supporting genres of another category. We present our approach to overcome this limitation and apply this approach to composing newspaper articles in the domain of fine arts.

# Composing Discourse based on Genre Characteristics

Kateryna Falkovych[1], Frank Nack[1,2]

Centrum voor Wiskunda en Informatica[1]

P.O.Box 94079

NL-1090 GB Amsterdam

the Netherlands

Firstname.Lastname@cwi.nl


V2_[2]

Institute for the Unstable Media

Eendrachtsstraat 10

3012 XL Rotterdam

The Netherlands

## Abstract

The availability of semantically enriched data in repositories of larger content providers offers a means for new ways of multimedia presentation authoring. Existing (semi-)automatic content composition environments explore limited numbers of genres and concentrate on template-based approaches for composing discourse structures for these genres. We analyze a number of genres based on identifiable genre characteristics. We show that existing template-based approaches to content composition support the essay and biography genres that belong to one characteristic genre category but fail in supporting genres of another category. We present our approach to overcome this limitation and apply this approach to composing newspaper articles in the domain of fine arts.

# 1 Introduction

Currently content providers, such as broadcasters[1] or museums[2], are moving towards enriching their content with semantic descriptions that allow offering larger freedom in interaction with the content of their multimedia repositories[3]. To enable this freedom they require approaches to content composition that allow users to become active authors of the multimedia content. Such approaches should provide efficient access to the content and encourage users in more frequent interactions with the media material. Responding to this demand a number of systems developed semantic-based approaches to support automatic or semi-automatic content composition scenarios.

Automatic approaches, such as [2], [12], [14] generate a multimedia presentation on request based on user preferences about a genre and a topic. Here decisions about multimedia content to be collected and specific structuring of the content are left to the system. Semi-automatic environments, such as [1], [9], [15], offer authoring support for the complete process of content composition where an author takes decisions about what content will be selected and how this content will be composed. The main emphasis of both authoring processes is to ensure the generation of a coherent presentation. An essential concept to achieve coherence in a presentation is genre. The approaches aim to support content composition for a large collection of genres in order to meet the needs of various authors.

Existing systems explore the biography and essay genres. They apply template-based approaches to guide the creation of discourse structures for these genres. To enable support of a larger number of genres we investigate a newspaper article genre. We analyze the differences between the genres based on the identifiable genre characteristics. The analysis suggests that essays and biographies essentailly differ from news articles. These essential differences cause unsuitability of the exiting template-based approaches for the composition of discourse structures for the newspaper article genres and other genres with similar characteristics. In this paper we go beyond the existing template-based approaches to overcome their limitations with respect to wider genre representation within (semi-)automatic multimedia presentation authoring.

The reminder of the paper is structured as follows. In Section 2 we outline a genre model, which allows us to describe recognizable genre characteristics important for the discourse composition tasks in multimedia authoring. Section 3

---

[1]http://www.bbc.co.uk/
[2]http://www.rijksmuseum.nl, http://www.metmuseum.org
[3]http://www.chip-project.org/, http://www.multimedian.nl

presents related work and our own experiences with template-based approaches for discourse structure composition. It also discusses limitations of these approaches for supporting various genres. Section 4 presents an alternative approach to discourse structure composition. The paper concludes in Section 5 with the evaluation of the proposed approach and directions for future work.

# 2 Genre - a rapprochement

To enable composition of discourse structures for various genres in an automatic or semi-automatic multimedia presentation generation process we need to have an understanding about the specific characteristics of each genre. These characteristics allow us to identify a genre and can serve as building blocks or as a reference during the discourse structure composition. The identification of recognizable genre characteristics requires a model in terms of which we can discuss and compare different genres. The following subsection presents a genre model. Subsection 2.2 discusses the differences between various genres in terms of this model.

## 2.1 Genre: content, form, function

The genre concept is applicable across different categories of artefacts, such as documents, novels, theater plays, films or multimedia presentations. It is believed that a genre represents established communication patterns that help readers or viewers to recognize and interpret information more efficiently [7]. Modern genre theory [6] sees genre as an evolving system. As a result there are no stable genre definitions and genre classifications. Yet, a genre is commonly defined as the combination of <content, form and function>:

- **Content** is interpreted as the subject of an artefact, events or characters contained in it. Content is often described in relation to form, so that content is a substance of an artefact while form covers its shape and structure.
- **Form** is usually seen as a visual appearance of an artefact, its structure, as manifested by its specific formatting and layout [19]. For media-rich artefacts form also includes specific design solutions (e.g. colors, relative sizes of different multimedia elements combined on one pages) and a mode of the presentation (e.g. an html page or a video).

- A particular combination of content and form identifies a specific **function** of the created artefact. For example, if we combine content (persons, computer science concepts, cities, countries, dates, companies) with form (author fields, title field, publishing date field, publisher field), we get the reference list genre which has the function of providing the background literature list and placing the topic in context. Genres can have other social (a newspaper article: to inform) or organizational (a dictionary: to collect the words of a language and their translations into another language) functions.

Content is crucially important for genre identification and is the core of the content composition issue. As shown in [19], the semantic content is the major constituent for genre recognition. Recognizable differences between genres that are important for our further analysis are based on specific characteristics of the content element. Therefore, we investigate this element in more detail.

Content element is a structure in itself that follows the same organisational principle as outlined for genre, namely its elements distinguish between content, form and function[4]. For clarity reasons we distinguish these elements as: <content$^c$, form$^c$ and function$^c$>.

Content$^c$ represents conceptual topics of the artefact. We distinguish between main topics, which form the subject of a presentation, and related topics, which provide background information, details or elaborate information. Thus, every related topic plays a particular role within the discourse flow, which can be regarded as the function$^c$ of a related topic with regard to the main topic. Besides, each topic is manifested in a particular way. For example, a topic can be conveyed as a paragraph of text and this text can be a quote. Additionally, the text inside a paragraph can be written in various writing styles (e.g. formal vs. informal language). Such aspects can be regarded as form$^c$ of a topic.

In a semantic-based systems the discourse structure composition process uses knowledge contained within a semantic graph. A semantic graph consists of domain concepts containing domain classes (Artist, Painting) and instances of these classes (e.g. Mondrian, Chrysanthemum) and semantic relations between them (e.g. Artist creates Painting). The mapping between the <content$^c$, form$^c$ and function$^c$>elements and elements of a semantic graph is shown below:

- **Content$^c$** is a collection of conceptual elements from which a discourse structure can be composed. These elements are expressed via domain concepts.

---

[4]A similar interpretation of the content element for the news genre is provided by [3]

- **Form$^c$** is a particular expression of a conceptual element. A single domain concept can be expressed by a variety of media items that differ in media types (e.g text, image, video, audio) and their discourse roles (quote, description, example).
- **Function$^c$**: conceptual elements can represent main or related topics of a discourse structure. In a particular discourse structure a conceptual element representing the related topic plays a certain role with regard to the concept representing the main topic. Domain concepts can provide background, elaboration or discussion for the main topic.

Based on the identified characteristics of genres we consider differences between genres in the following subsection.

## 2.2 What makes different genres differ

In this subsection we discuss recognizable characteristics of different genres relevant for multimedia presentations in the domain of fine arts, such as biography, essay and news article. We consider how these characteristics affect the interdependencies between <content$^c$, form$^c$ and function$^c$>elements.

The **biography** genre focuses on a list of facts such as birth, education, work, relationships and death, and establishes a complex insight of a personality including intimate details of experiences. The principal identification aspect of the biography genre is that it describes the life of a person. Hence, the genre is predominantly recognizable by its content$^c$. Form$^c$ of this content can vary. An author might use interviews, diaries or private letters. These materials can be presented in a textual form (the classic written biography) or in multimedia forms of biographies, where visual images add new dimensions in a presentation. With regard to the function$^c$ element the concepts used in a biography can alter their role in the discourse flow. For example, in a biography we can present information about education as background for professional activities. Alternatively, we can use personal relations as the background for describing the professional life if the former had an influence on it. But even though the form$^c$ and function$^c$ elements can change shapes, the important aspect of the biography is that particular concepts are mentioned, such as birth, upbringing, work, relationships and death, and those should be presented in sufficient detail.

The **essay** genre is a short piece that not only provides information about a particular topic but specifically aims at treating this topic from an author's personal point of view. An author usually selects a number of related topics to support his

arguments and opinions. As described by Aldous Huxley in the preface to his essay collection [13], an essay is an extreme variable genre that can be built on the three categories, namely the personal and the autobiographical; the objective, factual, concrete-particular; and the abstract-universal. This suggests that the essay is oriented towards offering a reader a description of the particular content while allowing for a free exploitation of form$^c$ and function$^c$ elements depending on its category. A personal viewpoint can influence functions of the related topics in an essay. A particular theme or category might require a certain form. Therefore, for an essay, as for a biography, the content$^c$ element is of the major importance.

The **newspaper article** genre is mainly characterized by its discourse flow. It is known for its pyramid structure, where the first part of an article gives concise information about an event and each of the following parts elaborates further on related topics mentioned in the first part. This elaboration can take different forms. After the description of the main event, the next part of the article usually gives the background about the event, the next parts provide the elaboration or discussion [21]. Thus, the related topics present in an article should provide specific functions (e.g. background, elaboration) with regard to the main topic to create the pyramid discourse flow of an article. Therefore, the newspaper article genre is characterized by its function$^c$ element. The content$^c$ element plays a secondary role here, since a newspaper article on any topic will have similar function elements. The form$^c$ element is also variable since styles of writing can differ depending on the overall topic of a newspaper and expectations of its audience.

The same analysis holds for the **scientific article** genre that has the very distinctive function$^c$ elements. A scientific article contains the recognizable discourse flow with identifiable elements, namely abstract, introduction, related work, presentation of an approach and results, discussion, future work. These elements follow a well-known general scheme where only slight variations are allowed.

As a conclusion we can say that genres can be distinguished as those that are recognized by their content, such as the biography or the essay, or those that are recognized by their function, such as the newspaper article or the scientific article. Then, strategies for discourse structure composition should be able to create discourse structures that incorporate recognizable characteristics of a genre they represent. In the following section we discuss existing strategies for composing content-oriented genres. We argue that template-based approaches they apply cannot be adapted for supporting function-oriented genres.

# 3 Strategies for various genres support

Current experiences with applying semantic web technology to multimedia authoring result in a number of strategies for creating biographies of artists [12], [14]. These strategies use a template-based approach for discourse structure composition. We applied a template approach to essay creation for a number of content composition tasks within a semi-automatic multimedia authoring environment [8], [9]. As the previous section showed, both of these genres are content-oriented. In the following subsection we describe the template-based approaches to composition of content-oriented discourse structures.

## 3.1 Creating content-oriented discourse structures

The **Artequakt** project [14] aims at automatically generating artist biographies from multimedia data extracted from the web and stored in a knowledge base. The multimedia data is annotated with the concepts from a dynamically populated domain ontology which contains the CIDOC [16] ontology. Discourse structures for biographies are presented using human-authored templates. A template consists of queries to the knowledge base that are composed using domain classes and relationships between them. Each query retrieves data about one aspect of an artist's life. Queries within a template are arranged in substructures that define variations of the ordering preferences. For instance, queries withing substructures can be called sequentially (Sequence), alternatively (Concept) or with regard to the set priorities (Level of Detail). An example of a template is shown below:

> **Sequence:**
> 1 d:Artist d:name Name
> 2 d:Artist d:place_of_birth Place
> 3 **Level of Detail:**
> priority=1 d:Artist d:influenced_by d:Person
> priority=2 d:Artist d:has_style d:Style
> *d:domain namespace*

**DISC** [12] uses the annotated repository of the Rijksmuseum Amsterdam [17] and a semantic graph encoded in RDF [23] to create multimedia presentations automatically on request. DISC, as Artequakt, explores the biography genre. The discourse structures for this genre are represented as dynamic rule-based templates. A template specifies a way of traversing the semantic graph. The dynamic nature of templates is achieved by creating them in a recursive manner.

A template is divided in a number of *narrative units* where each narrative unit represents a conceptual part of the discourse structure (e.g. personal data, private life, career). Within a narrative unit different domain classes can play roles of related characters. For instance, family members play roles of *disc:Spouse, disc:Son, disc:Father* related characters in the *disc:Private_Life* narrative unit. Besides, the template allows to specify for each related character whether it can play a role of the main character within certain narrative units. This type of recursion allows elaborating on a related character if the required data about this character can be found in the semantic graph. In the example template presented below the related character *disc:Spouse* can play a role of the main character in the *disc:Personal_Data* narrative unit:

> **disc:Personal_Data**
> d:Artist d:dateOfBirth Date
> d:Artist d:placeOfBirth Place
> d:Artist disc:role disc:MainCharacter
> **disc:Private_Life**
> d:Artist d:isMarried d:Person
> d:Person disc:role disc:Spouse
> disc:Spouse disc:role disc:MainCharacter ->applies to disc:Personal_Data
> *d:domain namespace*
> *disc:discourse namespace*

Thus, personal data of this character such as *d:date_of_birth* and *d:place_of_birth* will be found in the framework and included into the presentation. Such dynamic templates allow greater flexibility of otherwise predefined discourse structures since a number of related characters and available data about them will vary for each particular case.

**Samp*L*e** [9] is a semi-automatic multimedia presentation generation environment that aims at supporting authors during the complete multimedia presentation building process. The multimedia repository of the system covers fine arts. Samp*L*e uses a semantic graph that combines existing thesauri in the art domain (such as AAT [10] and ULAN [11] translated in OWL [24]) together with VRA schema [22] for annotating images and Dublin Core [5] as the top-level of the semantic graph. The semantic graph of Samp*L*e was extended with discourse role concepts describing the form aspect of multimedia material.

The process of presentation authoring is divided into four phases: topic identification, discourse structure building, media material collection and production

of the final-form presentation. During the first stage of Samp*L*e development[5] we created a support mechanism for one type of workflow where an author starts with defining a discourse structure and then has to collect media material to populate this structure. The support in this stage is oriented towards providing an author with a selection of discourse structures that suit her choice of topic and genre. In addition, the systems is able to suggest multimedia material appropriate for the chosen discourse structure based on semantic and discourse role annotations of media items:

> **Prologue** ->discourse roles: introduction, quote, definition
> 1. d:Movement
> **Main** ->discourse roles: description, elaboration, example
> 2. d:Style
> 3. d:Principle
> 4. ulan:Artist
> **Epilogue**->discourse roles: conclusion, quote
> 5. d:Movement
> *d: and ulan: are different domain namespaces*

During the second stage [8] we covered the inverse workflow where an author first collects media material from the repository while browsing and then the systems has to find a coherent structure to present selected material. The challenge here is to arrange the material within one of the discourse structures known to the system. Since foreseeing every possibility in mapping a varying set of media items to discourse structure is not a feasible task, the systems includes additional rules for resolving unforeseen situations.

Both of these approaches use templates to represent discourse structures. Discourse structures are created for the essay genre. The overall discourse flow is specified in a template using domain classes. The order of the classes ensures coherence of the discourse flow. The discourse structure extension mechanism allows to instantiate a template by posing queries to the semantic graph. In this way a number of related characters appearing in the discourse structure will vary according to the information found in the graph.

**Conclusions** The described systems use human-authored templates for discourse structures composition plus various additional mechanisms to allow dynamic changes within these structures. The main focus in creating discourse structures lies on the content$^c$ side. The discourse structure composition process

---

[5]http://www.cwi.nl/∼media/projects/CHIME/demos.html

9

is viewed from the content delivery perspective. This means that the process is mainly concerned with the problem what information should be presented in the discourse structure, and thus in the presentation. Such a view is appropriate for the essay and biography genres since these genres are oriented towards the presentation of a particular content and they do not impose strict requirements either on the roles of related topics (function$^c$) or on the particular way of expression (form$^c$). As discussed in section 2.2, the form and function elements within these genres change depending on the author, viewpoint or purpose.

This view is opposite to the situation when the discourse flow of a discourse structure is the main component defining a genre and thus is of the major importance for the discourse structure composition process. In this situation the function$^c$ element should guide the discourse structure composition. In the following subsection we discuss a problem of supporting composition of discourse structures for function-oriented genres. We show that content templates are not a feasible solution in this case.

## 3.2  Creating function-oriented discourse structures

We discuss the creation of function-oriented discourse structures using the newspaper article genre as a working example. We chose this genre, since it is a representative instance of function-oriented genres. It has a characteristic discourse flow, recognizable by its pyramid structure. The pyramid structure defines that the actual news event comes first followed by further details. Thus, this structure indicates to a reader which components of the article give essential information and which provide a sort of "additional reading". Components of this structure and their order can vary, some of them can be absent or merged. The recognition of these components might also depend on the formal or personal interpretation. Still the distinctions between the components in principle can be made based on general conventions  [20].

In an attempt to make these general conventions more tangible from a discourse analysis point of view, T. van Dijk  [21] proposed an analytical framework for the structures of news discourse. The schema of news consists of conventional categories that include: *Main Event(s), Context*, *History* and *Comments*. The conventional categories specify which function a topic of the content within each category should have. For instance, topics placed within the *Background* category should have such relationships with the main topic (the topic of the main event) that they provide background information for the main topic in the context of the current article. This suggests that even though the relationships between the

topics in the real world or concepts in a semantic graph do not change, our view on the function of each relation changes with respect to the context of the specific article. For example, the specific style a painter is using in his work can serve as *Background* in an article discussing his paintings. On the other hand, the information about this style can be used within the *Comments* category in an article about artists of a movement and their styles of work. Therefore, the set of topics forming the conceptual content of an article will change depending on the focus of an article and a specific instance of the main topic rather than its class. Consequently, we cannot create abstract templates that define which classes of topics (or classes of domain concepts) should be present within which structural elements of an article. If we use templates to approach the article composition problem, we would have to provide a template for each particular article we would like to create.

To address this problem we present an alternative approach to article composition. The approach adapts van Dijk's framework for news discourse to describe the discourse flow for an article applicable in a semantic-based context.

# 4 Alternative approach to content composition

We take van Dijk's framework of structural elements for an article as the basis for the composition process. As discussed in the previous section, the main topic of an article defines the article conceptually. To build a discourse structure for an article we have to propagate the discourse flow starting from the domain concept that represents *Main Event* (the main topic). On one hand the discourse should develop according to conventional categories of an article. On the other hand it should result in a coherent story.

## 4.1 From article schema to CACs

To map the conventional categories of an article to concepts within the semantic graph we use a simple theory of conversation that argues that at any point in a conversation, there are only a few general categories of follow-up statements that constitute a natural continuation rather than a topic shift. These categories are called "conversational/associational categories" (CACs) [18]. We use the interpretations of the categories adapted by [4], namely $\text{Context}^{CACs}$, $\text{Specifics}^{CACs}$, $\text{Analogy}^{CACs}$, $\text{Alternative}^{CACs}$.

- The *Context* category of the article schema has a direct relations to the Context$^{CACs}$ category.

- The *History* category can be seen as further elaboration on the *Main Event* or its *Context*. To provide information for the *History* category we can give more specific information about the main topic (Specifics$^{CACs}$) or provide information about similar occasions (Analogy$^{CACs}$).

- The *Comments* category provides additional details that can be expressed via Specifics$^{CACs}$ or Analogy$^{CACs}$ and a place for stating opinions. Alternative$^{CACs}$ is especially relevant for the *Comments* category, since it provides an alternative view on the situation, an alternative approach or result (see Figure 2 on page 28).

## 4.2   From CACs to semantic graph

The selected CACs provide a clearer view on how we can propagate a discourse flow starting with the main topic:

- The context of a concept is usually another concept that gives a more general perspective. Thus, we can express Context$^{CACs}$ by finding a domain concept related to the main topic that is more general than the main topic.

- In the same way, we can express Specifics$^{CACs}$ by finding a more specific related concept within the semantic graph.

- Analogy$^{CACs}$ can be expressed by finding two concepts of the same type A (belonging to the same class) that are both related to one concept of another class B. Class A should be more specific than class B.

- Two concepts are Alternative$^{CACs}$ to each other if: they belong to the same class A in the semantic graph; they are related correspondently to the other two concepts which belong to the same class B. Class A should be more specific than class B.

These examples suggest that if we can distinguish between general and specific concepts in the semantic graph then we can map the CACs to the domain concepts. In a semantic graph it is not possible to identify which of the concepts are more general or more specific. In contrast, in a tree parents can be more general concepts and children more specific. In order to make a semantic graph interpretable in tree terms, we introduce generality/specificity measurements of domain concepts. These measurements address the class level of the semantic graph.
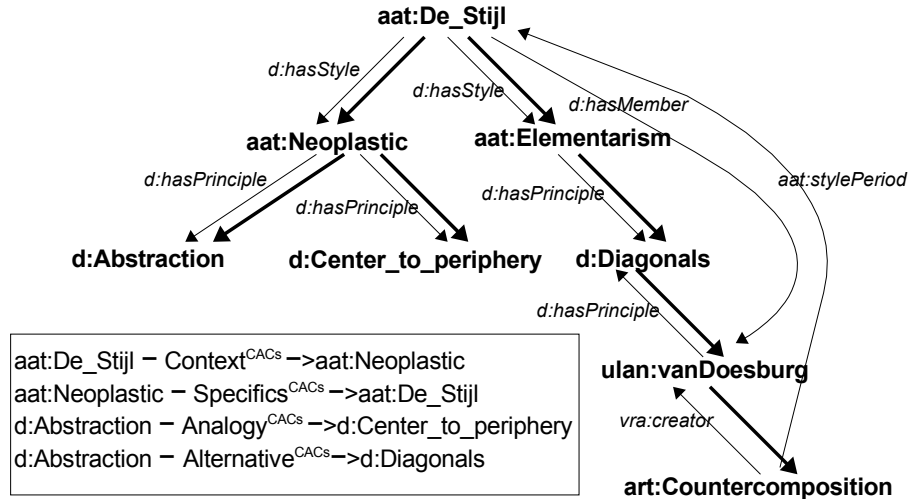
12

Figure 1: The mapping between the CACs and the semantic graph

They are domain-specific and should be provided by a domain expert. Within the art domain the movement class is the most general, since it unites other domain concepts: there are particular styles used within the movement, particular artists or artworks belong to the movement. Figure 1 shows a part of the semantic graph where thin lines represent semantic relationships between the concepts, bold lines represent tree-like relations. The relations between the domain concepts enclosed in the frame represent examples of Context$^{CACs}$, Specifics$^{CACs}$, Analogy$^{CACs}$ and Alternative$^{CACs}$.

## 4.3 Article composition process

For each of the selected CACs we created a rule with the corresponding name which can identify context concepts, specifics concepts, analogical and alternative concepts for a given concept using semantic relationships between the concepts in the semantic graph and generality/specificity measurements between them. Then the process of article composition takes the following shape. In order to find concepts appropriate to fill in *Context* category (context concepts), we apply *context rule* to the concept representing *Main Event* (the main topic) (see Figure 2 on page

28). For the *History* category we apply *specifics* rule to the main topic or analogy rule to the main topic or context concepts. The *Comments* category is filled in by applying *specifics, analogy* or *alternative rule* to the main topic or *analogy* and *alternative* to the context concepts. The rules are applied repeatedly since either of them conceptually matches the meaning of the corresponding conventional category as discussed earlier. Besides, it helps to not overrestrict the search space. Having more options for mapping conventional categories of an article to the domain concepts allows creating a number of discourse structures for an article on the specific topic.

In the general case each of these rules returns multiple results since, for example, there can be a number of concepts analogical (Analogy$^{CACs}$) to the current one in a semantic graph. Besides, application of multiple rules within each conventional category also contributes to the larger number of concepts retrieved. To choose the most appropriate concept with which to propagate the discourse flow we use coherence rules. A coherence rule analyzes the concepts that are already present in a discourse structure (*Current nodes*) and the set of concepts that are candidates for being included in the discourse structure (*Possible nodes*). After the analysis a set of *Next nodes* will be added to the discourse structure. We apply a set of the following coherence rules:

- *Repetition*: the repetition of concepts within a discourse structure is not allowed;

- *Consistency*: each following concept being added to a discourse structure has to have a semantic relation to the concept added at the previous step;

- *Pace*: all the concepts being added to a discourse structure should be within the scope of the main topic. This is achieved by identifying the highest general concept for the main topic using semantic relations and generality/specificity measurements of domain concepts. Only concepts that are directly (one semantic relations connects two concepts) or indirectly (two concepts are related via another concept) related to the highest general concept can be included into a discourse structure.

- *Succession*: do not include more specific concepts in a discourse structure if related to them more general concepts were not introduced. This rule gets applied only in cases where candidate concepts to be included in a discourse structure present too detailed information.

14

The whole process of article composition, the rules for expression CACs and coherence rules are implemented in Prolog. We currently integrate this functionality into the Samp*L*e [9] environment developed in our previous work.

# 5   Evaluation and conclusions

The described function-based process demonstrates the feasibility of our approach. The prototype engine is able to support the generation of discourse structures for the news article genre[6]. We use this genre as a representative example of the function-oriented category of genres. Although the presented discourse structure composition process contains the mapping which is specific for the news article genre, i.e. the mapping between the article schema and CACs, other mappings and developed rules are genre-independent and can be reused.

With regard to the proposed news article composition approach we can still see at least two major points for improvement. First, identification of a topic appropriate to be the main topic of an article is not trivial. In the process of generating a multimedia presentation on request an author in most cases is able to select a topic of the presentation out of a set of domain concepts present in a repository. For content-oriented genres almost any concept can serve as a main topic of the presentation. The main topic of an article is usually an event representing conflicting opinions or actions. Such an event is unique for each particular theme. Thus, the selection of the main topic of an article requires solid domain knowledge. Using the knowledge present in a semantic graph it is not possible to infer whether a concept can serve as a main topic of an article.

We can see two options for solving this problem. Either an author has to point out a relevant topic, or we need to include additional relations in a semantic graph to enable topic identification by the system. For instance, we could extend the set of domain relations with a "conflict" relation which will identify a pair of conflicting concepts. This could be conflicting principles of a movement or conflicting actions of artists involved in one movement. The system could search for such pairs and present an author with the choice of main topics for an article.

For the article composition approach this would mean that we need to include a strategy for adding the conflict pair of the main topic in a discourse structure. We could search for the most appropriate place for the conflict pair inside the created discourse structure. Alternatively we could reconsider the article composition

---

[6]http://www.cwi.nl/~katya/MMSJ/articles.html contains examples of generated articles

approach in a way that an article gets composed having a prerequisite that the main topic and its conflict pair have to be included into a discourse structure.

The second general improvement is related to the role of a topic. As we discussed in Section 3.2, the role which a specific topic has with regard to the main topic, changes depending on a particular article. This problem is currently not addressed and is the subject of future investigations. Our view is that providing a strategy for including a conflict pair in the process will help to identify the context of the current discourse structure (e.g. by identifying the semantic distance between the conflict pair concepts) and will allow to shape the overall process of article composition in general and identification of functions of related topics in particular.

In our discussion about various genres and their composition strategies we made a distinction between the different composition approaches for content- and function-oriented genres. However we on no account suggest that only content-oriented approaches can support content-oriented genres. Even though existing approaches are able to support the creation of content-oriented genres with content templates, it does not mean this is the only solution. Content-oriented genres can also be considered them from the function perspective. For example, in a biography the information about parents of a person can have an elaborative function if the life of this person was not influenced mainly by his parents. On the other hand, this information can have the function of a background for the other events in the biography, if the parents' influence was essential for the main character. Our assumption is that using function-oriented approaches for discourse structure building of content-oriented genres can contribute to the quality of discourse structures being created. A discourse structure could become more tailor-made for a particular main topic. To make such a function-oriented approach realizable we have to be able to identify (a) discourse flow(s) for content-oriented genres and to map a discourse flow to the domain concepts and relations. This is the direction of our ongoing research.

# 6   Acknowledgments

# References

[1] B. Barry. The Mindful Camera: Common Sense for Documentary Videography. In *Proceedings of the 11th ACM International Conference on Multimedia*, pages 648–649, November 2003.

[2] S. Bocconi, F. Nack, and L. Hardman. Supporting the Generation of Argument Structure within Video Sequences. In *Proceedings of the sixteenth ACM Conference on Hypertext and Hypermedia 2005*, pages 75–84, September 2005.

[3] L. Breure. Reuse of Content and Digital Genres. In H. van Oostendorp, L. Breure, and A. Dillon, editors, *Creation, Use and Deployment of Digital Information*, pages 27–53. Lawrence Erlbaum Associates, 2005.

[4] C. Cleary and R. Bareiss. Practical methods for automatically generating typed links. In *Proceedings of the the seventh ACM conference on Hypertext*, pages 31–41, Bethesda, Maryland, United States, March 1996.

[5] Dublin Core Community. Dublin Core Element Set, Version 1.1, 2003.

[6] D. Duff. *Modern Genre Theory*. Pearson Education Limited, Edinburgh Gate, United Kingdom, 2000.

[7] T. Erickson. Rhyme and Punishment: The Creation and Enforcement of Conventions in an On-Line Participatory Limerick Genre. In *Proceedings of the 32rd Hawaii International Conference on System Science*. IEEE Computer Society, January 5-8, 1999.

[8] K. Falkovych and S. Bocconi. Creating a Semantic-based Discourse Model for Hypermedia Presentations: (Un)discovered Problems. In *Workshop on Narrative, Musical, Cinematic and Gaming Hyperstructure*, Salzburg, Austria, September 2005.

[9] K. Falkovych and F. Nack. Context Aware Guidance for Multimedia Authoring: Harmonizing Domain and Discourse Knowledge. *Multimedia Systems Journal, Special issue on Multimedia System Technologies for Educational Tools, S. Acton, F. Kishino, R. Nakatsu, M. Rauterberg & J. Tang eds.*, 2006 to be published.

[10] Getty Research Institute. Art & Architecture Thesaurus (Online). http://www.getty.edu/research/tools/vocabulary/aat/, 2000. Version 2.0.

[11] Getty Research Institute. Union List of Artist Names (Online). http://www.getty.edu/research/conducting_research/vocabularies/ulan/, 2000. Version 2.0.

[12] J. Geurts, S. Bocconi, J. van Ossenbruggen, and L. Hardman. Towards Ontology-driven Discourse: From Semantic Graphs to Multimedia Presentations. In *Second International Semantic Web Conference (ISWC2003)*, pages 597–612, Sanibel Island, Florida, USA, October 20-23, 2003.

[13] A. Huxley. *Collected Essays*. Harper, New York, 1959.

[14] S. Kim, H. Alani, W. Hall, P. Lewis, D. Millard, N. Shadbolt, and M. Weal. Artequakt: Generating Tailored Biographies with Automatically Annotated Fragments from the Web. Presented at the Semantic Authoring, Annotation and Knowledge Markup (SAAKM) 2002 Workshop at the 15th European Conference on Artificial Intelligence (ECAI 2002), Lyon, France.

[15] F. Nack and W. Putz. Saying What It Means: Semi-Automated (News) Media Anotation. *Multimedia Tools and Applications*, 22(3):263 – 302, March 2004.

[16] N.Crofts, D.M.Dionissiadou, and M.Stiff. Definition of the cidoc object-oriented conceptual reference model. International Organization for Standardization, Technical report, 2000.

[17] Rijksmuseum Amsterdam. Rijksmuseum Amsterdam Website. http://www.rijksmuseum.nl.

[18] R. C. Schank. *Cognitive Science, Vol. 1*, chapter Rules and Topics in Conversation, pages 421–441. 1977.

[19] E. G. Toms and D. G. Campbell. Genre as Interface Metaphor: Exploiting Form and Function in Digital Environments. In *Thirty-Second Annual Hawaii International Conference on System Sciences-Volume 2*, page 2008, 1999.

[20] T. A. van Dijk. *News as discourse*. Lawrence Erlbaum Associates, Hillsdale, NJ, 1988.

[21] T. A. van Dijk. *Handbook of Qualitative Methods in Mass Communication Research*, chapter The interdisciplinary study of news as discourse, pages 108–120. Routledge, London, 1991.

[22] Visual Resources Association. Visual Resources Association Website.

[23] W3C. Resource Description Framework (RDF) Model and Syntax Specification. W3C Recommendations are available at http://www.w3.org/TR, February 22, 1999.

[24] W3C. Web Ontology Language (OWL) Reference Version 1.0. Work in progress. W3C Working Drafts are available at http://www.w3.org/TR, 12 November 2002.
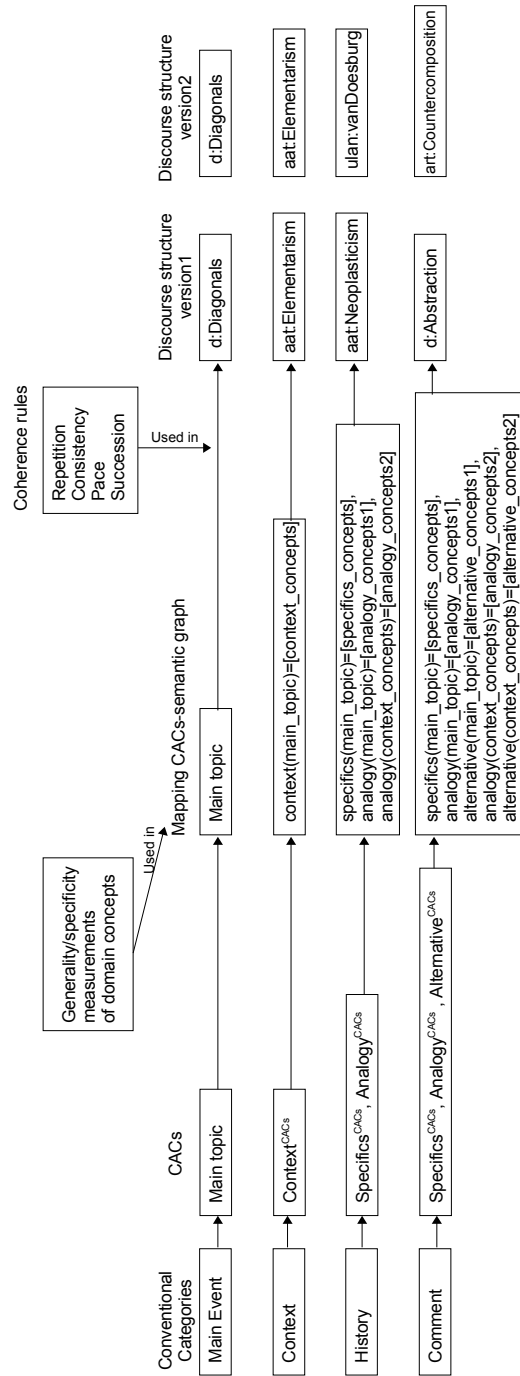
Figure 2: The process of article composition