Enhancement of Knowledge Resources and Discovery by Computation of Content Factors

Claus-Peter Rückemann Westfälische Wilhelms-Universität Münster (WWU), Leibniz Universität Hannover, North-German Supercomputing Alliance (HLRN), Germany Email: ruckema@uni-muenster.de

Abstract—This paper presents a methodology for data description and analysis, the Content Factor (CONTFACT). The Content Factor method can be applied to arbitrary data and content and it can be adopted for many purposes. Normed factors and variants can also support data analysis and knowledge discovery. This paper presents the algorithm, introduces into the norming of Content Factors, and discusses examples and a practical case study and implementation based on long-term knowledge resources, which are continuously in development. The methodology is used for advanced processing and also enables methods like data rhythm analysis and characterisation. It can be integrated with complementary methodology, e.g., classification and allows the application of advanced computing methods. The goal of this research is to create a general and flexible methodology for data description and analysis that can be used with huge structured and even unstructured data resources, allows an automation, and can therefore also be used for longterm multi-disciplinary knowledge.

Keywords–Data-centric Knowledge Processing; Content Factor (CONTFACT) method; Data Rhythm Analysis; Universal Decimal Classification; Advanced Computing.

I. INTRODUCTION

Information systems handling unstructured as well as structured information are lacking means for data description and analysis, which is data-centric and can be applied in flexible ways. In the late nineteen nineties, the concept of in-text documentation balancing has been introduced with the knowledge resources in the LX Project. Creating knowledge resources means creating, collecting, documenting, and analysing data and information. This can include digital objects, e.g., factual data, process information, and executable programs, as well as realia objects. Long-term means decades because knowledge is not isolated, neither in space nor time. All the more, knowledge does have a multi-disciplinary context.

Therefore, after integration knowledge should not disintegrate, instead it should be documented, preserved, and analysed in context. The extent increases with growing collections, which requires advanced processing and computing. Especially the complexity is a driving force, e.g., in depth, in width, and considering that parts of the content and context may be continuously in development. Therefore, the applied methods cannot be limited to certain algorithms and tools. Instead there are complementary sets of methods.

The methodology of computing factors [1] and patterns [2] being representative for a certain part of content was consid-

ered significant for knowledge resources and referred material. Fundamentally, a knowledge representation is surrogate. It enables an entity to determine consequences without forcing an action. For the development of these resources a definitionsupported, sortable documentation-code balancing was created and implemented.

The Content Factor (CONTFACT) method advances this concept and integrates a definition-supported sortable documentation-code balancing and a universal applicability. The Content Factor method is focussing on documentation and analysis The Content Factor can contain a digital 'construction plan' or a significant part of digital objects, like sequenced DeoxyriboNucleic Acid (DNA) does for biological objects [3]. Here, a construction plan is what is decided to be a significant sequence of elements, which may, e.g., be sorted or unsorted. Furthermore, high level methods, e.g., "rhythm matching", can be based on methods like the Content Factor.

Classification has proven to be a valuable tool for longterm and complex information management, e.g., for environmental information systems [4]. Conceptual knowledge is also a complement for data and content missing conceptual documentation, e.g., for data based on ontologies used with dynamical and autonomous systems [5].

Growing content resources means huge amounts of data, requirements for creating and further developing advanced services, and increasing the quality of data and services. With growing content resources content balancing and valuation is getting more and more important.

This paper is organised as follows. Section II summarises the state-of-the-art and motivation, Sections III and IV introduce the Content Factor method and an example for the application principle. Section V shows implemented Content Factor examples, explains flags, definition sets, and norming. Section VI provides the results from an implementation case study, showing complementary properties and complex scenarios. Section VII discusses aspects of processing and computation. Sections VIII and IX present and evaluation and main results, summarise the lessons learned, conclusions and future work.

II. STATE-OF-THE-ART AND MOTIVATION

Most content and context documentation and knowledge discovery efforts are based on data and knowledge entities. Knowledge is created from a subjective combination of different attainments, which are selected, compared and balanced against each other, which are transformed, interpreted, and used in reasoning, also to infer further knowledge. Therefore, not all the knowledge can be explicitly formalised.

Knowledge and content are multi- and inter-disciplinary long-term targets and values [6]. In practice, powerful and secure information technology can support knowledge-based works and values. Computing goes along with methodologies, technological means, and devices applicable for universal automatic manipulation and processing of data and information. Computing is a practical tool and has well defined purposes and goals.

Most measures, e.g., similarity, distance and vector measures, are only secondary means [7], which cannot cope with complex knowledge. Evaluation metrics are very limited, and so are the connections resulting from co-occurences in given texts, e.g., even with Natural Language Processing (NLP), or clustering results in granular text segments [8].

Evaluation can be based on word semantic relatedness, datasets and evaluation measures, e.g., the WordSimilarity 353 dataset (EN-WS353) for English texts [9]. The development of Big Data amounts and complexity up to this point show that processing power is not the sole solution [10]. Advanced longterm knowledge management and analytics are on the rise.

Value of data is an increasingly important issue, especially when long-term knowledge creation is required, e.g., knowledge loss due to departing personnel [11]. Current information models are not able to really quantify the value of information. Due to this fact one of the most important assets [12], the information, is often left out [13]. Today a full understanding of the value of information is lacking. For example, free Open Access contributions can bear much higher information values than contributions from commercial publishers or providers.

For numberless application scenarios the entities have to be documented, described, selected, analysed, and interpreted. Standard means like statistics and regular expression search methods are basic tools used for these purposes.

Anyhow, these means are not data-centric, they are volatile methods, delivering non-persistent attributes with minimal descriptive features. The basic methods only count, the result is a number. Numbers can be easily handled but in their solelity such means are quite limited in their descriptiveness and expressiveness.

Therefore, many data and information handling systems create numbers of individual tools, e.g., for creating abstracts, generating keywords, and computing statistics based on the data. Such means and their implementations are either very basic or they are very individual.

The pool of tools requires new and additional methods of more universal and data-centric character – for structured and unstructured data.

New methods should not be restricted to certain types of data objects or content and they should be flexibly usable in combination and integration with existing methods and generally applicable to existing knowledge resources and referenced data. New methods should allow an abstraction, e.g., for the choice of definitions as well as for defined items.

III. THE CONTENT FACTOR

The fundamental method of the Basic Content Factor (BCF), $\kappa_{\rm B}$ – "Kappa-B" –, and the Normed Basic Content Factor (NBCF), $\overline{\kappa}_{\rm B}$, can be described by simple mathematical notations. For any elements o_i in an object o, holds

$$o_i \in o. \tag{1}$$

The organisation of an object is not limited, e.g., a reference can be defined an element. For $\kappa_{\rm B}$ of an object *o*, with elements o_i and the count function *c*, holds

$$\kappa_{\rm B}(o_i) = c(o_i) \,. \tag{2}$$

For $\overline{\kappa}_{\rm B}$ of an object *o*, for all elements *n*, with the count function *c*, holds

$$\overline{\kappa}_{\mathrm{B}}(o_i) = \frac{c(o_i)}{\sum_{i=1}^{n} c(o_i)} \,. \tag{3}$$

All normed κ for the elements o_i of an object o sum up to 1 for each object:

$$\sum_{i=1}^{n} \overline{\kappa}_{\mathrm{B}}(o_i) = 1.$$
(4)

For a mathematical representation counting can be described by a set o and finding a result n, establishing a one to one correspondence of the set with the set of 'numbers' 1, 2, 3, ..., n. It can be shown by mathematical induction that no bijection can exist between 1, 2, 3, ..., n and 1, 2, 3, ..., munless n = m. A set can consist of subsets. The method can, e.g., be applied to disjoint subsets, too. It should be noted that counting can also be done using fuzzy sets [14].

IV. APPLICATION EXAMPLE

The methodology can be used with any object, independent if realia objects or digital objects. Nevertheless, for ease of understanding the examples presented here are mostly considering text and data processing. Elements can be any part of the content, e.g., equations, images, text strings, and words. In the following example, "letters" are used for demonstrating the application. Given is an object with the sample content of 10 elements:

For this example it is suggested that A and Z are relevant for documentation and analysis. The relevant elements, AAAZ, in an object of these 10 elements for element A means 3/10 normed so the full notation is

AAAZ/10 with
$$\overline{\kappa}_{\rm B}(A) = 3/10$$
 and $\overline{\kappa}_{\rm B}(Z) = 1/10$. (6)

In consequence, the summed value for AAAZ/10 is

$$\overline{\kappa}_{\rm B}(\mathsf{A},\mathsf{Z}) = 4/10. \tag{7}$$

AAAZ in an object of 20 elements, for element A means 3/20 normed, which shows that it is relatively less often in this object. 3/22 for element A for this object means this object or

an instance in a different development stage, e.g., at a different time or in a different element context. The notation

$$\{i_1\}, \{i_2\}, \{i_3\}, \dots, \{i_n\}/n$$
 (8)

of available elements holds the respective selection where $\{i_1\}, \{i_2\}, \{i_3\}, \ldots, \{i_n\}$ refers to the definitions of element groups. Elements can have the same labels respectively values. From this example it is easy to see that the method can be applied independent from a content structure.

V. CONTENT FACTOR EXAMPLES

The following examples (Figures 1, 2, 4, 3, 5) show valid notations of the Normed Basic Content Factor $\overline{\kappa}_{B}$, which were taken from the LX Foundation Scientific Resources [15]. The LX Project is a long-term multi-disciplinary project to create universal knowledge resources. Application components can be efficiently created to use the resources, e.g., from the Geo Exploration and Information (GEXI) project. Any kind of data can be integrated. Data is collected in original, authentic form, structure, and content but data can also be integrated in modified form. Creation and development are driven by multifold activities, e.g., by workgroups and campaigns. A major goal is to create data that can be used by workgroups for their required purposes without limiting long-term data to applications cases for a specific scenario. The usage includes a targetted documentation and analysis. For the workgroups, the Content Factor has shown to be beneficial with documentation and analysis. There are countless fields to use the method, which certainly depend on the requirements of the workgroups. For the majority of use cases, especially, selecting objects and comparing content have been focus applications. With these knowledge resources multi-disciplinary knowledge is documented over long time intervals. The resources are currently already developed for more than 25 years. A general and portable structure was used for the representation.

```
1 CONTFACT:20150101:MS: {A} {A} {G} {G} {G} /2900
```

```
2 CONTFACT:20150101:M:{A}:=Archaeology|Archeology
```

```
3 CONTFACT: 20150101:M: {G}:=Geophysics
```

Figure 1. NBCF $\bar{\kappa}_B$ for an object, core notation including the normed CONTFACT and definitions, braced style.

The Content Factor can hold the core, the definitions, and additional information. The core is the specification of $\kappa_{\rm B}$ or $\overline{\kappa}_{\rm B}$. Definitions are assignments used for the elements of objects, specified for use in the core.

Here, the core entry shows an International Standards Organisation (ISO) date or optional date-time code field, a flag, and the CONTFACT core. The definitions hold a date-time code field, flag, and CONTFACT definitions or definitions sets as shown here. Definition sets are groups of definitions for a certain Content Factor. The following examples show how the definition sets work.

```
CONTFACT:20150101:MS:AAG/89
```

```
CONTFACT: 20150101:M:A:=Archaeology | Archeology
```

```
CONTFACT:20150101:M:G:=Geophysics
```

1	CONTFACT: 20150101: MU: A{Geophysics}{Geology}/89
---	---

2 **CONTFACT**: 20150101:**M**:A:=Archaeology|Archeology

```
3 CONTFACT:20150101:M:{Geophysics}:=Geophysics|
Seismology|Volcanology
4 CONTFACT:20150101:M:{Geology}:=Geology|
```

```
Palaeontology
```

Figure 3. NBCF $\overline{\kappa}_{\rm B}$ for an object, core notation including the normed CONTFACT and definitions, mixed style.

```
1 CONTFACT:20150101:MU:{Archaeology}{Geophysics}/120
```

```
2 CONTFACT:20150101:M:Archaeology:=Archaeology|
Archeology
```

```
3 CONTFACT: 20150101:M: Geophysics:=Geophysics
```

Figure 4. NBCF $\overline{\kappa}_B$ for an object, core notation including the normed CONTFACT and definitions, multi-character non-braced style.

```
1 CONTFACT:20150101:MU:vvvvaSsC/70
```

- 2 CONTFACT: 20150101:M:v:=volcano
- 3 CONTFACT:20150101:M:a:=archaeology
- 4 CONTFACT:20150101:M:S:=Solfatara
- 5 CONTFACT: 20150101:M:s:=supervolcano
- 6 CONTFACT:20150101:M:C:=Flegrei

Figure 5. NBCF $\overline{\kappa}_{\rm B}$ for an object from a natural sciences collection, multi-case non-braced style.

Definitions can, e.g., be valid in braced, non-braced, and mixed style. Left values can have different labels, e.g., uppercase, lowercase, and mixed style can be valid. Figure 6 shows an example using Universal Decimal Classification (UDC) notation definitions.

```
1 CONTFACT:20150101:MS:{UDC:55}/210
2 CONTFACT:20150101:M:{UDC:55}:=Earth Sciences. Geological
    sciences
```

Figure 6. NBCF $\overline{\kappa}_{\rm B}$ for an object from a natural sciences collection, UDC notation definitions, braced style.

Conceptual knowledge like UDC can be considered in many ways, e.g., via classification and via description.

A. Flags

Content Factors can be associated with certain qualities. Sample flags, which are used with core, definition, and additional entries are given in Table I.

TABLE I. SAMPLE FLAGS USED WITH CONTFACT ENTRIES.		
Purpose	Flag	Meaning
Content Factor quality	U	Unsorted
	S	Sorted
Content Factor source	М	Manual
	А	Automated
	Н	Hybrid

The CONTFACT core entries can have various qualities, e.g., unsorted (U) or sorted (S). Unsorted means in the order in which they appear in the respective object. Sorted means in a different sort order, which may also be specified. CONTFACT entries can result from various workflows and procedures, e.g., they can be created on manual base (M) or on automated base

Figure 2. NBCF $\overline{\kappa}_B$ for an object, core notation including the normed CONTFACT and definitions, non-braced style.

(A). If nothing else is specified the flag refers to the way object entries were created. Content Factor quality refers to core entries, source also refers to the definitions and information.

The Content Factor method provides the specified instructions. The required features with an implementation can, e.g., implicitly require large numbers of comparisons, resulting in highly computationally intensive workflows on certain architectures. It is the choice of the user to weighten between the benefits and the computational efforts, and potentially to provide suitable environments.

B. Definition sets

Definition sets for object elements can be created and used very flexibly, e.g., word or string definitions. Therefore, a reasonable set of elements can be defined for the respective purpose, especially:

- Definition sets can contain appropriate material, e.g., text or classification.
- Groups of elements can be created.
- Contributing elements can be subsummarised.
- Definition sets can be kept persistent and volatile.
- Definition set elements can be weighted, e.g., by parameterisation of context-sensitive code growth.
- Context sensitive definition sets can be referenced with data objects.
- Content can be described with multiple, complementary definition sets.
- Any part of the content can be defined as elements.

The Content Factors can be computed for any object, e.g., for text and other parts of content. Nevertheless, the above definition sets for normed factors are intended to be used with one type of elements.

C. Normed application

 $\overline{\kappa}_{\rm B}$ is a normed quantity. Norming is a mathematical procedure, by which the interesting quantity (e.g., vector, operator, function) is modified by multiplication in a way that after the norming the application of respective functionals delivers 1. The respective $\overline{\kappa}_{\rm B}$ Content Factor can be used to create a weighting on objects, e.g., multiplying the number of elements with the respective factor value.

VI. IMPLEMENTATION

The implementation has been created for the primary use with knowledge resources' objects (lxcontfact). This means handling of any related content, e.g., documentation, keywords, classification, transliterations, and references. The respective objects were addressed as Content Factor Object (CFO) (standard file extension .cfo) and the definition sets as Content Factor Definition (CFD) (standard file extension .cfd).

A. Case study: Computing complementation and properties

The following sequence of short examples shows a knowledge resources object (Figure 7), and three pairs of complementary CONTFACT definition sets and the according $\overline{\kappa}_{\rm B}$ computed for the knowledge resources object and respective definition sets (Figures 8 and 9; 10 and 11; 12 and 13).

1 object A %-GP%-XX%: object A [A, B, C, D, O]: 2 %-GP%-EN%: A B C D O 3 %-GP%-EN%: A B C D O 4 %-GP%-EN%: A B C D O 5 %-GP%-EN%: A B C D O 6 %-GP%-EN%: A B C D O	
--	--

Figure 7. Artificial knowledge resources object (LX Resources, excerpt).

The right parts are entry and keywords. Here, the algorithm can count in object entry name (right "object A"), keywords (in brackets), and object documentation (lower right block).

```
1 % (c) LX-Project, 2015, 2016
2 {A}:=\bA\b
3 {0}:=\bO\b
```

Figure 8. CONTFACT definition set 1 of 3 (LX Resources, excerpt).

The definition set defines $\{A\}$ and $\{O\}$. The definitions are case sensitive for this discovery. We can compute $\overline{\kappa}_B$ (Figure 9) according to the knowledge resources object and definition set.

1	CONTFACT: BEGIN
2	CONTFACT: 20160117-175904: AU: {A} {A} {O} {A}
3	CONTFACT: 20160117-175904:AS: {A} {A} {A} {A} {A} {A} {O} {0} {0} {0} {0} {0} {32
4	CONTFACT:20160117-175904:M:{A}:=\bA\b
5	CONTFACT:20160117-175904:M:{0}:=\b0\b
6	CONTFACT: 20160117-175904:M:STAT: OBJECTELEMENTSDEF=2
7	CONTFACT: 20160117-175904:M:STAT: OBJECTELEMENTSALL=32
8	CONTFACT: 20160117-175904:M:STAT:OBJECTELEMENTSMAT=13
9	CONTFACT: 20160117-175904:M:STAT: OBJECTELEMENTSCFO=.40625000
10	CONTFACT: 20160117-175904:M:STAT: OBJECTELEMENTSKWO=2
11	CONTFACT: 20160117-175904:M:STAT: OBJECTELEMENTSLAN=1
12	CONTFACT:20160117-175904:M:INFO:OBJECTELEMENTSOBJ=object A
13	CONTFACT:20160117-175904:M:INFO:OBJECTELEMENTSDCM=(c) LX-Project, 2015, 2016
14	CONTFACT: 20160117-175904:M: INFO: OBJECTELEMENTSMTX=LX Foundation Scientific
	Resources; Object Collection
15	CONTFACT: 20160117-175904:M: INFO: OBJECTELEMENTSAUT=Claus-Peter R\"uckemann
16	CONTFACT: END

Figure 9. NBCF $\overline{\kappa}_{\rm B}$ computed for knowledge resources object and definition set 1 (LX Resources, excerpt).

The result is shown in a line-oriented representation, each line carrying the respective date-time code for all the core, statistics, and additional information. The second complementary set (Figure 11) defines $\{B\}$ and $\{D\}$.

- 1 % (c) LX-Project, 2015, 2016
- 2 {B}:=\bB\b
 3 {D}:=\bD\b
- 3 (U) =: {U}

Figure 10. CONTFACT definition set 2 of 3 (LX Resources, excerpt).

1	CONTFACT: BEGIN
2	CONTFACT: 20160117-175904:AU: {B} {D} {D} {D} {D} {D} {D} {D} {D} {D} {D
3	CONTFACT: 20160117-175904:AS: {B} {B} {B} {B} {B} {D} {D} {D} {D} {D} {D} {D} {A3
4	CONTFACT: 20160117-175904:M: {B}:=\bB\b
5	CONTFACT: 20160117-175904:M: {D}:=\bD\b
6	CONTFACT: 20160117-175904: M: STAT: OBJECTELEMENTSDEF=2
7	CONTFACT: 20160117-175904: M: STAT: OBJECTELEMENTSALL=33
8	CONTFACT: 20160117-175904: M: STAT: OBJECTELEMENTSMAT=12
9	CONTFACT: 20160117-175904: M: STAT: OBJECTELEMENTSCFO=. 37500000
10	
	Figure 11. NBCF $\overline{\kappa}_{\rm B}$ computed for knowledge resources object and
	definition set 2 (LX Resources, excerpt).
TI	he resulting $\overline{\kappa}_{\rm B}$ is shown in the excerpt (Figure 12).
11	ie resulting <i>n</i> _B is shown in the except (Figure 12).

\[\{\} (c) LX-Project, 2015, 2016

2 {C}:=\bC\b

Figure 12. CONTFACT definition set 3 of 3 (LX Resources, excerpt).

The third complementary set (Figure 13) defines $\{C\}$.

1	CONTFACT : BEGIN
2	CONTFACT:20160117-175905:AU:{C}{C}{C}{C}{C}/32
3	CONTFACT:20160117-175905:AS:{C}{C}{C}{C}{C}/32
4	CONTFACT:20160117-175905:M:{C}:=\bC\b
5	CONTFACT: 20160117-175905:M:STAT: OBJECTELEMENTSDEF=1
6	CONTFACT: 20160117-175905:M: STAT: OBJECTELEMENTSALL=32
7	CONTFACT: 20160117-175905:M:STAT:OBJECTELEMENTSMAT=6
8	CONTFACT:20160117-175905:M:STAT:OBJECTELEMENTSCFO=.18750000
9	

Figure 13. NBCF $\overline{\kappa}_B$ computed for knowledge resources object and definition set 3 (LX Resources, excerpt).

The sum of all elements considered for $\overline{\kappa}_B$ by the respective CONTFACT algorithm in an object is 100 percent. Here, the overall number of

- definitions is 2 + 2 + 1 = 5,
- elements is 32,
- matches is 13 + 12 + 6 = 31.

The sum of aggregated $\overline{\kappa}_{\rm B}$ values for all relevant elements results in 0.40625000 + 0.3750000 + 0.18750000 + 1/32 = 1.

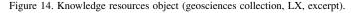
B. Case study: Complex resources and discovery scenario

The data used here is based on the content and context from the knowledge resources, provided by the LX Foundation Scientific Resources [15]. The LX knowledge resources' structure and the classification references [16] based on UDC [17] are essential means for the processing workflows and evaluation of the knowledge objects and containers.

Both provide strong multi-disciplinary and multi-lingual support. For this part of the research all small unsorted excerpts of the knowledge resources objects only refer to main UDC-based classes, which for this part of the publication are taken from the Multilingual Universal Decimal Classification Summary (UDCC Publication No. 088) [18] released by the UDC Consortium under the Creative Commons Attribution Share Alike 3.0 license [19] (first release 2009, subsequent update 2012).

The excerpts (Figures 14, 15, 16), show a CFO from the knowledge resources a CFD and the computed CONTFACT.

1 2 3	Vesuvius [Volcanology, Geology, Archaeology]: (lat.) Mons Vesuvius. (ital.) Vesuvio.
4	
	Volcano, Gulf of Naples, Italy.
5	Complex volcano (compound volcano). Stratovolcano, large cone (Gran Cono).
6	
7	The most well known antique settlements at the Vesuvius are
	Pompeji}, \lxidx{Herculaneum}, and \lxidx{Stabiae}.
8	s. also seismology, phlegra, Solfatara
9	%%IML: keyword: volcano, Vesuvius, Campi Flegrei, phlegra, scene of
,	
	fire, Pompeji, Herculaneum, volcanic ash, lapilli, catastrophe,
	climatology, eruption, lava, gas ejection, Carbon Dioxide
10	<pre>%%IML: UDC:[911.2+55]:[57+930.85]:[902]"63"(4+37+23+24)=12=14</pre>
11	
12	Object: Volcanic material.
13	Object-Type: Realia object.
14	Object-Location: Vesuvius, Italy.
15	Object-FindDate: 2013-10-00
16	Object-Discoverer: Birgit Gersbeck-Schierholz, Hannover, Germany.
17	Object-Photo: Claus-Peter Rückemann, Minden, Germany.
18	%%IML: media: YES 20131000 {LXC:DETAILM-} { UDC: (0.034)(044)770}
10	LXDATASTORAGE://img 3824.jpg
19	%%IML: UDC-Object:[551.21+55]:[911.2](37+4+23)=12
20	<pre>%%IML: UDC: 551.21 :: Vulcanicity. Vulcanism. Volcanoes. Eruptive</pre>
	phenomena. Eruptions
21	<pre>%%IML: UDC: 55 :: Earth Sciences. Geological sciences</pre>
22	<pre>%%IML: UDC: 911.2 :: Physical geography</pre>



Labels, language fields, and spaces were stripped. A knowledge object can contain any items required, e.g., including storing data, documentation, classification, keywords, algorithms, references, implementations, in any languages and representations, allowing support tables and algorithms. Examples of application scenarios for the Content Factor method range from libraries, natural sciences and archaeology, statics, architecture, risk coverage, technology to material sciences [20].

```
1 % (c) LX-Project, 2009, 2015
2 {Ve}:=Vesuvius
3 {Vo}:=\b[Vv]olcano
4 {Po}:=Pompe[ji]i
5 {UDC:55}:=Geology
6 {UDC:volcano}:=UDC.*\b911\b.*\b55\b
```

Figure 15. CONTFACT definition set (geosciences collection, LX, excerpt).

The definition sets can contain anything required for the definitions and additional information for the respective Content Factor implementation, e.g., definitions of elements and groups as well as comments. The left side defines the element used in the Content Factor and the right side states the matching element components. Left value and right value are separated by ":=" for an active definition.

1	CONTFACT: BEGIN
2	CONTFACT: 20160130-235804: AU: {Ve} {Vo} {UDC: 55: geology} {Ve} {Ve} {Vo} {Vo} {Vo} {Vo} {Vo} {Vo} {Vo} {Vo
	}{Vo}{Vo}{Ve}{Ve}{Po}{Ve}{Po}{Ve}{Ve}{Ve}{Vo}{Vo}{Vo}{Vo}{Vo}{UDC:volcano}{Vo}{
	Vo}/319
3	CONTFACT:20160130-235804:AS: {Po} {Po} {UDC:55:geology} {UDC:volcano} {Ve} {Ve} {Ve} {Ve}
	Vo}/319
4	CONTFACT:20160130-235804:M:{Ve}:=Vesuvius
5	CONTFACT:20160130-235804:M:{Vo}:=\b[Vv]olcano
6	CONTFACT:20160130-235804:M:{Po}:=Pompe[ji]i
7	CONTFACT:20160130-235804:M:{UDC:55:geology}:=Geology
8	CONTFACT:20160130-235804:M:{UDC:volcano}:=UDC.*\b911\b.*\b55\b
9	CONTFACT: 20160130-235804: M: STAT: OBJECTELEMENTSDEF=5
10	CONTFACT: 20160130-235804: M: STAT: OBJECTELEMENTSALL=319
11	CONTFACT: 20160130-235804: M: STAT: OBJECTELEMENTSMAT=28
12	CONTFACT: 20160130-235804: M: STAT: OBJECTELEMENTSCFO=.09180304
13	CONTFACT:20160130-235804:M:INFO:OBJECTELEMENTSDCM=(c) LX-Project, 2009, 2015
14	
15	CONTFACT: END

Figure 16. NBCF $\overline{\kappa}_{\rm B}$ computed for knowledge resources object and definition set (geosciences collection, LX Resources, excerpt).

The left value can include braces (e.g., curly brackets) in order to support the specification and identification of the left value. The right value can include common representations of pattern specification. The result of which can be seen from the computed CONTFACT.

The example patterns follow the widely used Perl (Practical Extraction and Report Language) regular expressions [21], e.g., \bar{b} for word boundaries and [...] and multiple choices of characters at a certain position.

C. Case study: Rhythm matching and core sequences

As soon as Content Factors have been computed for an object the patterns can be compared with pattern of other objects. The Content Factor method allows to compare the occurrences of relevant elements in objects in many ways. The following example shows the "rhythm matching" method for two computed unsorted CONTFACT core sequences (Figures 18, 19) for an object and a definition set (Figure 17).

1	\$ (c) LX-Project, 2009, 2015, 2016
2	{Am}:=\b[Aa]mphora
3	{Ce}:=[Cc]eramic
4	{Gr}:=\b[Gg]reek\b
5	{Pi}:=[Pp]itho[is]
6	{Ro}:=\b[Rr]oman\b
7	{Tr}:=[Tt]ransport
8	{Va}:=[Vv]ases
1	

Figure 17. Example of CONTFACT definition set, geoscientific and archaeological resources (LX Resources, excerpt).

CONTFACT:20160101-215751:AU: {Am} {Gr} {Ce} {Ro} {Am} {Gr} {Am} { Am} {Va} {Pi} {Tr}/474

Figure 18. CONTFACT rhythm matching: Computed core for same object (before modification) and definition set (LX Resources, excerpt).

```
\begin{array}{l} \textbf{CONTFACT: } 20160101 - 231806: \textbf{AU:} \left\{ \text{Am} \right\} \left\{ \text{Ce} \right\} \left\{ \text{Ro} \right\} \left\{ \text{Am} \right\} \left\{ \text{Am} \right\} \left\{ \text{Va} \right\} \left\{ \text{Tr} \right\} \left\{ \text{Ce} \right\} \left\{ \text{Tr} \right\} \left\{ \text{Ce} \right\} \left\{ \text{Tr} \right\} \left\{ \text{Ce} \right\} \left\{ \text{Am} \right\} \left\{ \text{Am} \right\} 488 \end{array}
```

Figure 19. CONTFACT rhythm matching: Computed core for same object (after modification) and definition set (LX Resources, excerpt).

The comparison shows that relevant passages were appended to the object (italics font). Relevant regarding the rhythm matching means relevant from the object and definition set. Even short sequences like $\{Am\}\{Gr\}\{Ce\}$ and even when sorted like $\{Am\}\{Ce\}\{Gr\}$ can be relevant and significant in order to compute factors and identify and compare objects. The Content Factor method does not have built-in or intrinsic limitations specifying certain ways of further use, e.g., with comparisons and analysis.

Unsorted CONTFACT are more likely to describe objects and quality, including their internal organisation. Sorted CON-TFACT tend to describe objects by their quantities, with reduced focus on their internal organisation.

Objects with larger amount of documentation maybe candidates for unsorted CONTFACT. Objects, e.g., with factual, formalised content maybe candidates for sorted CONTFACT. Combining several methods in a workflow is possible.

Anyhow, the further use of the CONTFACT core, e.g., sorting the core data for a certain comparison, is a matter of application and purpose with respective data.

VII. PROCESSING AND COMPUTATION

A. Scalability, modularisation, and dynamical use

The algorithms can be used for single objects as well as for large collections and containers, containing millions of entries each. The computation routines allow a modularised and dynamical use.

The parts required for an implementation computing a Content Factor can be modularised, which means that not only a Content Factor computation can be implemented as a module but even core, definitions, and additional parts can be computed by separate modules.

A sequence of routine calls used for examples in this case study shows the principle and modular application of respective functions (Figure 20).

The modules create an entity for the implemented Content Factor (begin to end). They include labels, date, unsorted elements and so on as well as statistics and additional information.

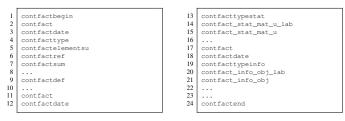


Figure 20. Sequence of modular CONTFACT routines for lxcontfact implementation (LX Resources, excerpt).

Application scenarios may allow to compute Content Factors for many objects in parallel. Content Factors can be computed dynamically as well as in batch mode or "pre-computed". Content Factors can be kept volatile as well as persistent. Everything can be considered a set, e.g., an object, a collection, and a container. Therefore, an implementation can scale from single on the fly objects to millions of objects, which may also associated with pre-computed Content Factors.

B. Parallelisation and persistence

There is a number of modules supporting computation based on persistent data, e.g., in collections and containers. The architecture allows task parallel implementations for multiple instances as well as highly parallel implementations for core routines. Examples are collection and container decollators, collection and container slicers, collection and container atomisers, formatting modules, and computing modules for (intermediate) result matrix requests.

Content Factor data can easily be kept persistent and dynamically. The algorithms and workflows allow the flexible organisation of data locality, e.g., central locations and with compute units, e.g., in groups or containers.

VIII. EVALUATION

The case study has shown that the formal description can be implemented very flexibly and successful (lxcontfact). Content Factors can be computed for any type of data. The Content Factor is not limited to text processing or even NLP, term-frequencies, and statistics. It has been successfully used with long term knowledge resources and with unstructured and dynamical data. The Content Factor method can describe arbitrary data in a unique form and supports data analysis and knowledge discovery in many ways, e.g., complex data comparison and tracking of relevant changes.

Definition sets can support various use cases. Examples were given from handling single characters to string elements. Definitions can be kept with the Content Factor, together with additional Content Factor data, e.g., statistics and documentation. Any of this Content Factor information has been successfully used to analyse data objects from different sources. The computation of Content Factors is non invasive, the results can be created dynamically and persistent. Content Factors can be automatically computed for elements and groups of large data resources. The integration with data and knowledge resources can be kept non invasive to least invasive, depending on the desired purposes. Knowledge objects, e.g., in collections and

containers, can carry and refer to complementary information and knowledge, especially Content Factor information, which can be integrated with workflows, e.g., for discovery processes.

The benefits and usability may depend on the field of application and the individual goals. The evaluation refers to the case context presented, which allows a wide range of freedom and flexibility. The benefits for the knowledge resources are additional means for documentation of objects. In detail, the benefits for the example workflows were improved data-mining pipelines, due to additional features for comparisons of objects, integrating developing knowledge resources, and creating and developing knowledge resources. In practice, the computation of Content Factors has revealed significant benefits for the creation and analysis of large numbers of objects and for the flexibility and available features for building workflows, e.g., when based on long-term knowledge objects. In addition, creators, authors, and users of knowledge and content have additional means to express their views and valuation of objects and groups of objects. From the computational point of view, the computation of Content Factors can help minimise the recurrent computing demands for data.

IX. CONCLUSION

This paper introduced a methodology for data description and analysis, the Content Factor (CONTFACT) method. The paper presents the formal description and examples, a successful implementation, and a practical case study. It has been shown that the Content Factor is data-centric and can describe and analyse arbitrary data and content, structured and unstructured. Data-centricity is even emphasized due to the fact that the Content Factor can be seamlessly integrated with the data. The data locality is most flexible and allows an efficient use of different computing, storage, and communication architectures.

The method can be adopted for many purposes. The Content Factor method has been successfully applied for knowledge processing and analysis with long-term knowledge resources, for knowledge discovery, and with variable data for system operation analysis. It enables to specify a wide range of precision and fuzziness for data description and analysis and also enables methods like data rhythm analysis and characterisation, can be integrated with complementary methodologies, e.g., classifications, concordances, and references. Therefore, the method allows weighting data regarding significance, promoting the value of data. The method supports the use of advanced computing methods for computation and analysis with the implementation. The computation and processing can be automated and used with huge and even unstructured data resources. The methodology allows an integrated use with complementary methodologies, e.g., with conceptual knowledge like UDC. It will be interesting to see various Content Factor implementations for individual applications, e.g., dynamical classification and concordances. Future work concentrates on advanced analysis and automation for different application scenarios, e.g., object comparisons, multi-lingual discovery, and concordance discovery.

ACKNOWLEDGEMENTS

We are grateful to the "Knowledge in Motion" (KiM) long-term project, Unabhängiges Deutsches Institut für Multidisziplinäre Forschung (DIMF), for partially funding this implementation, case study, and publication and to its senior scientific members, especially to Dr. Friedrich Hülsmann, Gottfried Wilhelm Leibniz Bibliothek (GWLB) Hannover, to Dipl.-Biol. Birgit Gersbeck-Schierholz, Leibniz Universität Hannover, and to Dipl.-Ing. Martin Hofmeister, Hannover, for fruitful discussion, inspiration, practical multi-disciplinary case studies, and the analysis of advanced concepts. We are grateful to all national and international partners in the GEXI cooperations for their constructive and trans-disciplinary support.

REFERENCES

- C.-P. Rückemann, "Advanced Content Balancing and Valuation: The Content Factor (CONTFACT)," Knowledge in Motion Longterm Project, Unabhängiges Deutsches Institut für Multidisziplinäre Forschung (DIMF), Germany; Westfälische Wilhelms-Universität Münster, Münster, 2009, Project Technical Report.
- [2] C.-P. Rückemann, "CONTCODE A Code for Balancing Content," Knowledge in Motion Long-term Project, Unabhängiges Deutsches Institut für Multidisziplinäre Forschung (DIMF), Germany; Westfälische Wilhelms-Universität Münster, Münster, 2009, Project Technical Report.
- [3] F. Hülsmann and C.-P. Rückemann, "Content and Factor in Practice: Revealing the Content-DNA," KiM Summit, October 26, 2015, Knowledge in Motion, Hannover, Germany, 2015, Project Meeting Report.
- [4] C.-P. Rückemann, "Integrated Computational and Conceptual Solutions for Complex Environmental Information Management," in The Fifth Symposium on Advanced Computation and Information in Natural and Applied Sciences, Proceedings of The 13th International Conference of Numerical Analysis and Applied Mathematics (ICNAAM), September 23–29, 2015, Rhodes, Greece, Proceedings of the American Institute of Physics (AIP). AIP Press, 2015, ISSN: 0094-243X, (in press).
- [5] D. T. Meridou, U. Inden, C.-P. Rückemann, C. Z. Patrikakis, D.-T. I. Kaklamani, and I. S. Venieris, "Ontology-based, Multi-agent Support of Production Management," in The Fifth Symposium on Advanced Computation and Information in Natural and Applied Sciences, Proceedings of The 13th International Conference of Numerical Analysis and Applied Mathematics (ICNAAM), September 23–29, 2015, Rhodes, Greece, Proceedings of the American Institute of Physics (AIP). AIP Press, 2015, ISSN: 0094-243X, (in press).
- [6] C.-P. Rückemann, F. Hülsmann, B. Gersbeck-Schierholz, P. Skurowski, and M. Staniszewski, Knowledge and Computing. Post-Summit Results, Delegates' Summit: Best Practice and Definitions of Knowledge and Computing, September 23, 2015, The Fifth Symposium on Advanced Computation and Information in Natural and Applied Sciences, The 13th International Conference of Numerical Analysis and Applied Mathematics (ICNAAM), September 23–29, 2015, Rhodes, Greece, 2015.
- [7] O. Lipsky and E. Porat, "Approximated Pattern Matching with the L_1, L_2 and L_∞ Metrics," in 15th International Symposium on String Processing and Information Retrieval (SPIRE 2008), November 10–12, 2008, Melbourne, Australia, ser. Lecture Notes in Computer Science (LNCS), vol. 5280. Springer, Berlin, Heidelberg, 2008, pp. 212–223, Amir, A. and Turpin, A. and Moffat, A. (eds.), ISSN: 0302-9743, ISBN: 978-3-540-89096-6, LCCN: 2008938187.
- [8] G. Ercan and I. Cicekli, "Lexical Cohesion Based Topic Modeling for Summarization," in Proceedings of The 9th International Conference on Computational Linguistics and Intelligent Text Processing (CICLing 2008), February 17–23, 2008, Haifa, Israel, ser. Lecture Notes in Computer Science (LNCS), vol. 4919. Springer, Berlin, Heidelberg, 2008,

pp. 582–592, Gelbukh, A. (ed.), ISSN: 0302-9743, ISBN: 978-3-540-78134-9, LCCN: 2008920439, URL: http://link.springer.com/chapter/ 10.1007/978-3-540-78135-6_50 [accessed: 2016-01-10].

- [9] G. Szarvas, T. Zesch, and I. Gurevych, "Combining Heterogeneous Knowledge Resources," in Proceedings of The 12th International Conference on Computational Linguistics and Intelligent Text Processing (CICLing 2011), February 20–26, 2011, Tokyo, Japan, ser. Lecture Notes in Computer Science (LNCS), vol. 6608 and 6609. Springer, Berlin, Heidelberg, 2011, pp. 289–303, Gelbukh, A. (ed.), ISSN: 0302-9743, ISBN: 978-3-642-19399-6, DOI: 10.1007/978-3-642-19400-9, LCCN: 2011921814, URL: http://link.springer.com/chapter/ 10.1007/978-3-642-19400-9_23 [accessed: 2016-01-10].
- [10] A. Woodie, "Is 2016 the Beginning of the End for Big Data?" Datanami, 2016, January 5, 2016, URL: http://www.datanami.com/2016/01/05/is-2016-the-beginning-of-the-end-for-big-data/ [accessed: 2016-01-10].
- [11] M. E. Jennex, "A Proposed Method for Assessing Knowledge Loss Risk with Departing Personnel," VINE: The Journal of Information and Knowledge Management Systems, vol. 44, no. 2, 2014, pp. 185–209, ISSN: 0305-5728.
- [12] R. Leming, "Why is information the elephant asset? An answer to this question and a strategy for information asset management," Business Information Review, vol. 32, no. 4, 2015, pp. 212–219, ISSN: 0266-3821 (print), ISSN: 1741-6450 (online), DOI: 10.1177/0266382115616301.
- [13] "The (Unknown) Value of Information (in German: Der (unbekannte) Wert von Information)," library essentials, LE_Informationsdienst, Dez. 2015 / Jan. 2016, 2015, pp. 10–14, ISSN: 2194-0126, URL: http://www. libess.de [accessed: 2016-01-10].
- [14] B. Kosko, "Counting with Fuzzy Sets," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. PAMI-8, no. 4, Jul. 1986, pp. 556–557, ISSN: 0162-8828, DOI: 10.1109/TPAMI.1986.4767822.
- [15] "LX-Project," 2016, URL: http://www.user.uni-hannover.de/cpr/x/ rprojs/en/#LX (Information) [accessed: 2016-01-01].
- [16] C.-P. Rückemann, "Enabling Dynamical Use of Integrated Systems and Scientific Supercomputing Resources for Archaeological Information Systems," in Proc. INFOCOMP 2012, Oct. 21–26, 2012, Venice, Italy, 2012, pp. 36–41, ISBN: 978-1-61208-226-4.
- [17] "UDC Online," 2015, URL: http://www.udc-hub.com/ [accessed: 2016-01-01].
- [18] "Multilingual Universal Decimal Classification Summary," 2012, UDC Consortium, 2012, Web resource, v. 1.1. The Hague: UDC Consortium (UDCC Publication No. 088), URL: http://www.udcc.org/udcsummary/ php/index.php [accessed: 2016-01-01].
- [19] "Creative Commons Attribution Share Alike 3.0 license," 2012, URL: http://creativecommons.org/licenses/by-sa/3.0/ [accessed: 2016-01-01].
- [20] F. Hülsmann, C.-P. Rückemann, M. Hofmeister, M. Lorenzen, O. Lau, and M. Tasche, "Application Scenarios for the Content Factor Method in Libraries, Natural Sciences and Archaeology, Statics, Architecture, Risk Coverage, Technology, and Material Sciences," KiM Strategy Summit, March 17, 2016, Knowledge in Motion, Hannover, Germany, 2016.
- [21] "The Perl Programming Language," 2016, URL: https://www.perl.org/ [accessed: 2016-01-10].