

Computing Optimised Result Matrices for the Processing of Objects from Knowledge Resources

Claus-Peter Ruckemann

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

Abstract—The aim of this paper is to discuss and summarise the main results on computing optimised result matrices from the practical creation of long-term multi-disciplinary and multi-lingual knowledge resources. Structuring big data is the essential process, which has to precede creating and implementing algorithms. The knowledge resources implement structure and features and can be integrated most flexibly into information and computing system components. Main elements are so called knowledge objects, which can consist of any content and context documentation and can employ a multitude of means for description and referencing of objects used with computational workflows. Core attributes are a faceted universal classification and various content views and attributes. Developing workflow implementations for various purposes requires to compute result matrices from the objects and referred knowledge, e.g., from geosciences, archaeology, physics, and information technology. The purposes can require individual processing means, complex algorithms, and a base of big data collections. Advanced discovery workflows can easily demand large computational requirements for High End Computing (HEC) resources supporting an efficient implementation. This paper presents some major methodologies and statistics instruments, which have been developed and successfully integrated. The combination of instruments and resources allows to flexibly compute optimised result matrices for discovery processes in information systems, expert and decision making components, search engine algorithms, and fosters the further development of the long-term knowledge resources.

Keywords—*Knowledge Processing; Result Matrix; Optimisation; Computing; Statistics; Classification; UDC; Big Data; High End Computing; Knowledge Resources; Knowledge Discovery.*

I. INTRODUCTION

Knowledge resources are the basic components in complex integrated systems. Their target is mostly to create a long-term multi-disciplinary knowledge base for various purposes. Request and selection processes result in requirements for computing result matrices from the available information and data. Optimisation in the context of result matrices means “improved for a certain purpose”. Here, the certain purpose is given by the target and intention of the application scenario, e.g., requests on search results or associations. Therefore, improving the result matrices is a very multi-fold process and “optimising result matrices” primarily refers to the content and context but in second order also to the workflows and algorithms. The major means presented here contributing to the optimisation are classification and statistics, based on the knowledge resources. The employed knowledge resources can

provide any knowledge documentation and additional information on objects and knowledge references, e.g., from natural sciences and decision making. Any data used in case studies is embedded into millions of multi-disciplinary objects, including dynamical and spatial information and data files.

It is necessary to develop logical structures in order to govern the existing unstructured and structured big data today and in future, especially in volume, variability, and velocity and to keep the information addressable on long-term. Preparing and structuring big data is the essential process, which has to precede creating and implementing algorithms. The systematic, methodological, and “clean” big data knowledge preparation and structuring must generally be named as largest achievement in this context and can be considered by far the most significant overall contribution [1]. The creation and optimisation of respective algorithms is of secondary importance, the more the data must be considered for long-term knowledge creation as, e.g., the benefits of most of those implementations depend on a certain generation of computing and storage architectures, which change all few 4–6 years.

Workflows based on these objects and facilities have been created for different applications. The knowledge resources can make sustainable and vital use of Object Carousels [2] in order to create knowledge object references and modularise the required algorithms [3]. This provides a universal means for improving coverage, e.g., dark data, and quality within the workflow. Secondary resources being available for data, information, and knowledge integration, besides Integrated Information and Computing System (IICS) applications, allow for workflows and intelligent components on High End Computing (HEC) and High Performance Computing (HPC) resources [4], [5]. This paper presents the up-to-date experiences with selected components for structures and workflows.

This paper is organised as follows. Section II introduces the previous work with methodologies and components used, Sections III and IV present the implemented means and statistics fundamentals integrated. Sections V, VI, and VII discuss the implementation environment and evaluate main results, and summarise the lessons learned, conclusions and future work.

II. METHODOLOGIES AND COMPONENTS EMPLOYED

The data used here is based on the content and context from the knowledge resources, provided by the LX Foundation Scientific Resources [6], [7]. The LX structure and the classification references based on UDC [8], [9] are essential means

for the processing workflows and evaluation of the knowledge objects and containers. Both provide strong multi-disciplinary and multi-lingual support.

An instructive example for an archaeological and geoscientific use case, deploying knowledge resources, classification, references, and Object Carousels has been recently published [2]. With this research the presentation complements the use case by an important methodology, statistics for intermediate result matrices, usable in any associated workflow. In order to get an overview, the following practical example for a specific workflow as part of an application component shows how result matrices for requests can be computed iteratively.

- 1) Application component request,
- 2) Object search (i.e., knowledge objects, classification, references, associations),
- 3) Creation of intermediate result matrices,
- 4) Iterative and alternating matrix element creation (i.e., based on intermediate result matrices, object search, referenced content, classification, and statistics),
- 5) Creation of result matrix,
- 6) Application component response.

The workflow will mostly be linear if the used algorithms are linear and the data involved is fixed in number and content.

The knowledge objects are under continuous development for more than twenty-five years. The classification information has been added in order to describe the objects with the ongoing research and in order to enable more detailed documentation in a multi-disciplinary and multi-lingual context.

Classification is state-of-the-art with the development of the knowledge resources, which implicitly means that the classification is not created statically or even fixed. It can be used and dynamically modified on the fly, e.g., when required by a discovery workflow description. Representations and references can be handled dynamically with the context of a discovery process. So, the classification can be dynamically modelled with the workflow context. The applied workflows and processing are based on the data and extended features developed for the Gottfried Wilhelm Leibniz resources [10].

Mathematical statistics is a central means for data analysis [11], [12]. It can be of huge benefits when analysing regularities and patterns when used for machine learning with information system components [13]. It is a valuable means deployed in natural sciences and has been integrated in multi-disciplinary humanities-based disciplines, e.g., in archaeology [14]. The span of fields for statistics is not only very broad but statistics itself goes far beyond a simple “tool” status [15].

Methodological means, which have been created in order to be deployed for regular use are workflows improving result quantity and result quality, various filters, universal classifications, statistics applications, manually documented resources’ components, integration interfaces for knowledge resources, comparative methods, combination of several means. The methodologies with the knowledge resources are based on computational methods, processing, classification and structuring of multi-disciplinary knowledge, systematic documentation, long-term knowledge creation, vitality of data concepts, sustainable resources architecture, and collaboration frameworks. In the past, many algorithms have been developed

and implemented [6], [7] for supporting different targets, e.g., silken criteria, statistics, classification, references and citation evaluation, translation, transliteration, and correction support, regular expression based applications, phonetic analysis support, acronym expansions, data and application assignments, request iteration, centralised and distributed discovery, and automated and manual contributions to the workflow.

A. Structure and classification

The key issues for computing result matrices from knowledge resources are that they require long-term tasks on efficiently structuring and classifying content and context. The classification, which has shown up being most important with complex multi-disciplinary long-term classification with practical simple and advanced applications of knowledge resources is the Universal Decimal Classification (UDC) [16]. According to Wikipedia currently about 150,000 institutions, mostly libraries and institutions handling large amounts of data and information, e.g., the ETH Library (Eidgenössische Technische Hochschule), are using basic UDC classification worldwide [17], e.g., with documentation of their resources, library content, bibliographic purposes on publications and references, for digital and realia objects. Just regarding the library applications UDC is present in more than 144,000 institutions and 130 countries [18]. Further operational areas are author-side content classifications and museum collections.

UDC allows an efficient and effective processing of knowledge data. UDC provides facilities to obtain a universal and systematical view on the classified objects. UDC in combination with statistical methods can be used for analysing knowledge data for many purposes and in a multitude of ways. With the knowledge resources in this research handling 70,000 classes, for 100,000 objects and several millions of referenced data then simple workflows can be linear but the more complex the algorithms get the workflows will mostly become non-linear. They allow interactive use, dynamical communication, computing, decision support, and pre- and postprocessing, e.g., visualisation. The classification deployed for documentation [19] is able to document any object with any relation, structure, and level of detail as well as intelligently selected nearby hits and references. Objects include any media, textual documents, illustrations, photos, maps, videos, sound recordings, as well as realia, physical objects, such as museum objects. UDC is a suitable background classification, for example: The objects use preliminary classifications for multi-disciplinary content. Standardised operations used with UDC are coordination and addition (“+”), consecutive extension (“/”), relation (“:”), order-fixing (“::”), subgrouping (“[]”), non-UDC notation (“*”), alphabetic extension (“A-Z”), besides place, time, nationality, language, form, and characteristics.

B. Statistics implementation for the knowledge resources

A vast range of statistics, e.g., mathematical statistics, can be deployed based on the knowledge resources. The application of mathematical statistics benefits from an increased number of probes or elements. Probes can result from measurements,

e.g., from applied natural sciences and from available material. In many cases, without further analysis a distribution or result may seem random. If the accumulation of an occurrence may indicate a regularity or a rule then this may correlate with a statistical method. Many cases require that statistical results have to be verified for realness. This can be done checking against experience and understanding and using mathematical means, e.g., computing probabilities based on probes.

Statistics have been used for steering the development of the resources. Classification and keyword statistics support the optimisation of the quality of data within the knowledge resources. Counts of terms, references, homophones, synonyms and many more support the improvement of the discovery workflows. Comparisons of content with different language representations increase the intermediate associated result matrices for a discovery process. The created knowledge resources' architecture is very flexible and efficient because the components allow a natural integration of multi-disciplinary knowledge. The processes of optimising a result matrix differ from a statistical optimisation by the fact that statistics is only one of the factors within the workflows.

III. IMPLEMENTED KNOWLEDGE RESOURCES' MEANS

The goals for the combination of statistics and classification are, for example:

- Creating and improving result matrices.
- Decision making within workflows.
- Further development of knowledge resources.
- Extrapolation and prediction.

The implementation for the required flexible workflow creation and levels is shown in the following sketch (Figure 1).

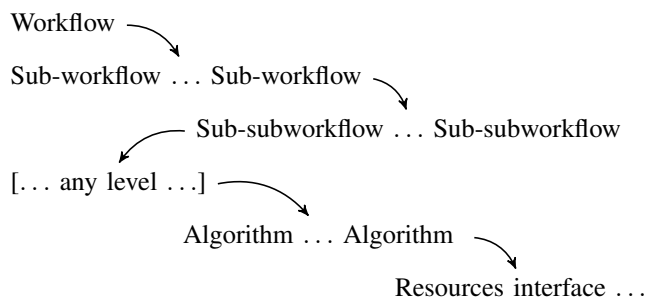


Figure 1. Workflow-algorithm sketch of the implementation.

The architecture is non-hierarchical. Any workflows can be applied in chains. Each workflow can use sub-workflows, these can use sub-subworkflows and so on. Each workflow can call or implement algorithms, e.g., for discovery processes, evaluation, and statistics. The workflows and algorithms can use or implement interfaces to the resources. The ellipses indicate that any step can be called or executed in parallel on HEC resources, e.g., in data-parallel or task-parallel processes, in any number of required instances.

An example for this is a “multi-probe parallelised optimisation” workflow, which generates an intermediate result matrix and uses the elements in order to create additional

results, all of which are combined for an overall optimised result matrix. The intermediate result matrices are deploying statistical, numerical methods, and various algorithms on base of additional knowledge and information resources.

As implementations of statistics are based on counting and numbers the statistics sub-workflows can deploy everything, e.g., any feature or attributes, which can be counted. Sources and means of statistics and computation are:

- Dynamical statistics on the internal and external content and context (e.g., overall statistics, keyword-, categories-, classification-, and media-statistics).
- Mathematics and formula on statistics from the content.
- Elements' statistics (structuring, content, references).
- Statistics based on UDC classification.
- UDC-based statistics computed from comparisons and associations of UDC groups and descriptions.
- Statistics based on any combination of classification, keywords, content, references, context, and computation.

Workflows based on the statistics can be type “semi-manually” or “automated”. Besides the major processing and optimisation goals descriptive statistics can be done with each workflow or sub-workflow. Any change of the means supported within a workflow can contribute to the optimisation of the result matrix. Suitable and appropriate means have to be determined for best supporting the goals of the respective step in the workflow. The implementation considers measuring the optimisation by quantity and quality of attributes and features, on intelligence-based and learning processes. With either use there is no general quality measure. Possible quality measures depend on purpose, view, and deployed means. In addition, the decision on these measures can be well supported by statistics, e.g., comparing result matrices from different workflows on the same request. Learning systems components can be used for capturing the success of different measures. The knowledge resources can contain equations and formulary of any grade of complexity. Due to the very high complexity level of the multi-disciplinary components it is necessary to use the basic instances for a comparison in this context of matrix statistics.

The following passages show basic excerpts of statistics objects (L^AT_EX representation) being part of the implemented knowledge resources. These statistics methods/equations are selected and shown mainly for two reasons: The selected methods are taken from the knowledge objects contained in the resources. These methods are used for result matrix calculations and compared with the evaluation in this research.

IV. STATISTICS: FUNDAMENTALS AND APPLICATION

Statistics on itself can rarely give an overall decisive answer on a question. Statistic means merely can be used as tools for supporting valuations and decisions. Statistics, probability, and distributions are valuable auxiliaries within workflows and integrated application components, e.g., on numbers of objects, spatial or georeferences, phonetic variations, and series of measurement values. Probability and statistics measures are used with integrated applications, e.g., with search requests, with seismic components (e.g., Median and Mean Stacks), which can also be implemented on base of the resources.

A. Basic algorithms applied with knowledge resources

The mean value, arithmetic mean or average M for n values is given by

$$M = \frac{1}{n} \sum_{\nu=1}^n x_{\nu} \quad (1)$$

Calculating the mean value is described by a linear operation. The median value or central value is the middle value in a size-depending sort order of a number of values. For making a statement on the extent of a group of values, the variance (“scattering”) can be calculated, with the mean deviation m and the squared mean deviation m^2 .

$$m^2 = \frac{1}{n} \sum_{\nu=1}^n (x_{\nu} - M)^2 = \overline{(x - M)^2} \quad (2)$$

For any value this holds $m^2(A) = m^2 + (M - A)^2$ When applying statistics, especially when calculating the propagated error, the following definition of the variance is used:

$$m^2 = \frac{1}{n-1} \sum_{\nu=1}^n (x_{\nu} - M)^2 \quad (3)$$

The mean deviation $\zeta(A)$ is defined as:

$$\zeta(A) = \overline{|x - A|} \text{ for which holds } \zeta(A) = \min. \text{ for } A = Z \quad (4)$$

The probable deviation or probable error ρ with the probable limits Q_1 and Q_3 is defined as:

$$\rho = \frac{Q_3 - Q_1}{2} \quad (5)$$

The relative frequency h_i is defined as:

$$h_i = \frac{n_i}{n}, \text{ then it holds } \sum_{i=1}^k h_i = 1 \quad (6)$$

where n_i is the class frequency, which means the number of elements in a class of which the middle element is x_i .

B. Distributions deployed with knowledge resources

A continuous summation results in the cumulative frequency distribution

$$H_i = \sum_{j=1}^i h_j \quad (7)$$

which gives the relative number for which holds $x \leq x_i$. H_i is a function discretely increasing from 0 to 1. The presentation results in a summation line. With steady variables, for which at an interval width of Δx the quotient $h_i/\Delta x$ nears a limit, one can calculate a frequency density $h(x_i)$ and for the summation frequency $H(x)$:

$$h(x_i) = \lim_{\Delta x \rightarrow 0} \frac{h_i}{\Delta x} \quad \text{and} \quad \frac{dH(x)}{dx} = h(x) \quad (8)$$

With statistical distributions the Gaussian normal distribution is of basic importance.

$$h(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad (9)$$

$H(x)$ can not be given “closed”. It can be shown that

$$K = \int_{-\infty}^{+\infty} h(x)dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{1}{2}x^2} dx = 1 \quad (10)$$

The Binominal distribution $w_k(s)$ is defined by

$$w_k(s) = \binom{k}{s} p^s q^{k-s} \quad (11)$$

The sum of the two binominal coefficients is equal to $\binom{k+1}{s}$. This is described by Pascals’ Triangle. It holds:

$$M = \sum_{s=0}^k w_k(s) \cdot s = kp \quad \text{and} \quad m = \sqrt{kqp} \quad (12)$$

Accordingly, the mean error of the mean value decreases proportional to $1/\sqrt{k}$. This describes the error propagation law.

$$h(X) = \frac{1}{\sqrt{2\pi m}} e^{-\frac{1}{2} \left(\frac{X-M}{m}\right)^2} \text{ for } -\infty < X < +\infty \quad (13)$$

From this Gaussian curves, binominal distributions, correlation coefficients and advanced measures can be developed.

C. Application of fundamental theorems of probability

The probability p is defined by:

$$p(x_i) = \lim_{n \rightarrow \infty} h_i = \lim_{n \rightarrow \infty} \frac{n_i}{n} \quad (14)$$

The classical definition of p_{classic} is:

$$p_{\text{classic}} = \frac{\text{number of favoured cases}}{\text{number of possible cases}} \quad (15)$$

The following can be said if independence is supposed. \vee means the logical OR, \wedge the logical AND.

Either-OR: If E_1, E_2, \dots, E_m are events excluding each other and the respective probabilities are p_1, p_2, \dots, p_m , then the probability for *either* E_1 *OR* E_2 *OR* \dots *OR* E_m is:

$$p(E_1 \vee E_2 \vee \dots \vee E_m) = p_1 + p_2 + \dots + p_m \quad (16)$$

As-well-as: If E_1, E_2, \dots, E_m are event pairwise independent from each other then the probability of E_1 *as well as* E_2 *as well as* \dots *as well as* E_m is:

$$p(E_1 \wedge E_2 \wedge \dots \wedge E_m) = p_1 \cdot p_2 \cdot \dots \cdot p_m \quad (17)$$

V. IMPLEMENTATION FOR THE RESULT MATRIX CASES

A. Measures for optimisation and purposes

The measures for optimisation are on the one hand object of the services and workflows but on the other hand they can be of concern for the knowledge resources themselves.

Conforming with the goals, measures for optimisation mean fitness for a purpose, e.g., search for a regularity with statistics and result matrices. After a search for regularities any statistical procedure benefits from checking against experiences and associating the procedure and result with a meaning. In many cases, e.g., “relevance” means numbers, uniqueness, proximity for objects, content, and attributes, e.g., terms.

Optimisation can be achieved by various means, e.g., by intelligent selection, by self-learning based optimisation, and by comparisons and statistics. The first measures include manual procedures and essences of results being stored for learning processes. They can also deploy comparisons and statistics, which also mean probability and distributions. This case study is focussed on comparisons and statistics applied with the knowledge resources. The subject of the statistics deals with the collection, description, presentation, and interpretation of data. Especially, the methodology can be based on computing more than the minimal number of comparisons, computing more than the minimal number of distributions, computing result matrices considering the mean of several distributions or extreme distributions. In the case of “relevance”, information on weighting may come from sources of different qualities.

The general steps with the knowledge resources, including external sources, can be summarised as: Knowledge resource requests, integrating search engine results (e.g., Google), integrating results more or less randomly, without explicit considerate classification and correlation between content and request, comparing the content of search result matrix elements with the knowledge resources result matrix containing classified elements, statistics on an accumulation of terms, selecting accumulated terms, elimination of less concentrated results, selecting the appropriate number of search results.

B. Sources and Structure: Knowledge resources

The full content, structure, and classification of the knowledge resources have been used. In the context of the case discussed here, the sources, which have been integrated and referenced with the knowledge resources consist of:

- Classical natural sciences data sources.
- Environmental and climatological information.
- Geological and volcanological information.
- Natural and man-made factor/event information.
- Data sets and compilations from natural sciences.
- Archaeological and historical information.
- Archive objects references to realia objects.
- Photo and video objects.
- Dynamical and non-dynamical computation of content.

The sources consist of primary and secondary data and are used for workflows, as far as content or references are accessible and policies, licenses, and data security do not restrict.

C. Classification and statistics in this sample case

Table I shows a small excerpt of resulting main UDC classification references practically used for the statistics with the knowledge resources in the example case presented here.

TABLE I. UNIVERSAL DECIMAL CLASSIFICATION OF STATISTICS FEATURES WITH THE KNOWLEDGE RESOURCES (EXCERPT).

UDC Code	Description
UDC:3	Social Sciences
UDC:310	Demography. Sociology. Statistics
UDC:311	Statistics as a science. Statistical theory
UDC:311.1	Fundamentals, bases of statistics
UDC:311.21	Statistical research
UDC:311.3	General organization of statistics. Official statistics
UDC:5	Mathematics. Natural sciences
UDC:519.2	Probability. Mathematical Statistics
UDC:531.19	Statistical mechanics
UDC:570.087.1	Biometry. Statistical study and treatment of biological data
UDC:615.036	Clinical results. Statistics etc.

The 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) [8] released by the UDC Consortium under the Creative Commons Attribution Share Alike 3.0 license [20] (first release 2009, subsequent update 2012).

As with any object the statistics features can be combined for facets and views for any classification subject. On the other hand statistics objects from the resources can be selected and applied. The listing (Figure 2) shows an excerpt intermediate object result matrix on statistics content.

1	ANOVA	[Statistics, ...]:
2		Analysis of Variance.
3	BIWS	[Whaling]:
4		Bureau of International Whaling Statistic.
5	GSP	[Geophysics]:
6		Geophysical Statistics Project.
7	Median	[Statistics]:
8		In the middle line.
9		s. also Median-Stack
10	Median-Stack	[Seismics]:
11		Stacking based on the median value of adjacent traces.
12	MSWD	[Mathematics]:
13		Mean Square Weighted Deviation.
14	MSA	[Abbreviation, GIS]:
15		Metropolitan Statistical Area.
16	MOS	[Abbreviation]:
17		Model Output Statistics.
18	MCDM	[GIS, GDI, Statistics, ...]:
19		Multi-Criteria Decision Making.
20	SHIPS	[Meteorology]:
21		Statistical Hurricane Intensity Prediction Scheme.
22	SAND	[Abbreviation]:
23		Statistical Analysis of Natural resource Data, Norway.

Figure 2. Intermediate object result matrix on “statistics” content.

Learning from this: The classifications used for this intermediate matrix are based on contributions from more than one discipline. The elements themselves do not necessarily have to contain a requested term because the classification

contributes. Several steps may be necessary in order to improve the matrix, e.g., selecting disciplines, time intervals on the entries, references, and associations. Because different content carries different attributes and features the evaluation can be used in comparative as well as in complementary context.

The implemented knowledge resources means of statistics and computation described above are integrated in the workflows, including classification, dating, and localisation of objects. In addition, probability distributions, linear and non-linear modelling, and other supportive tools are used within the workflow components.

D. Resulting numbers on processing and computing

The processing and computational demands per workflow instance result from the implementation scenarios. The following comparison (Table II) results from a minimal workflow request for a result matrix compared to a workflow request for a result matrix supporting classification views referring to UDC, supporting references and statistics on intermediate results. Both scenarios are based on the same number of elements and entries and can be considered atomic instances in a larger workflow. Views and result matrices can be created manually and automated in interactive and batch operation.

TABLE II. PROCESSING AND COMPUTATIONAL DEMANDS: 2 SCENARIOS, BASED ON 50000 OBJECT ELEMENTS AND 10 RESULT MATRIX ENTRIES.

<i>Scenario Workflow Request for Result Matrix</i>	<i>Value</i>
“geosciences archaeology” (minimal)	
Number of elements	50,000
Number of result matrix entries (defined)	10
Number of workflow operations	15
Wall time on one core	14 s
“geosciences archaeology” (UDC, references, statistics)	
Number of elements	50,000
Number of result matrix entries (defined)	10
Number of workflow operations	6,500
Wall time on one core	6,700 s

As the discussed scenarios are instances this means workflows based on n of these instances will at least require n -times the time for an execution on the same system. It must be remembered that the parallelisation will have a significant effect when workflows are created based on many of these instances when required in parallel. Without modifying the algorithms of the instances, which mostly means simplifying, the positive parallelisation effect for the workflows can be nearly linear. Besides the large requirements per instance with most workflows there are significant beneficial effects from parallelising even within single instances as soon as the number of comparable tasks based on the instances increases. A typical case where parallelisation within a workflow is favourable is the implementation of an application creating result matrices and being used with many parallel instances, e.g., with providing services. The number of 70,000 elementary UDC classes currently results in 3 million basic elements when only considering multi-lingual entries – without any combinations. With most isolated resources only several thousand combinations are used in practice each. The variety

and statistics are mostly deployed for decision processing, increasing quantity, and increasing quality. Many of the above cases require to compute more than one data-workflow set to create a decision. A review and an auditing process are mandatory for mission critical applications. The computational requirements can increase drastically with the computation of multiple workflows. Each workflow will consist of one or more processes, which can contain different configurations and parameters. Therefore, creating a base for an improved result matrix starts with creating several intermediate result matrices. With a ten process workflow, e.g., the possible configurations and parameters can easily lead to computing a reasonable set of thousands to millions of intermediate result matrices.

The objects and methods used can be long-term documented as knowledge objects. Nevertheless, there is explicitly no demand for a certain programming language. Even multiple implementations can be done with any object. The workflows and algorithms with the cases discussed here have been implemented as objects in Fortran, Perl, and Shell. Anyhow, the implementation of algorithms is explicitly not part of any core resources. It is the task of anyone having an application to do this and to decide on the appropriate means and methods.

VI. CASE RESULTS AND EVALUATION

Computing result matrices is an arbitrary complex task, which can depend on various factors. Applying statistics and classification to knowledge resources has successfully provided excellent solutions, which can be used for optimising result matrices in context of natural sciences, e.g., geosciences, archaeology, volcanology or with spatial disciplines, as well as for universal knowledge. The method and application types used for optimisation imply some general characteristics when putting discovery workflows into practice regarding components like terms, media, and other context (Table III).

TABLE III. RESULTING PER-INSTANCE-CALLS FOR METHOD AND APPLICATION TYPES ON OPTIMISATION WITH KNOWLEDGE DISCOVERY.

<i>Type</i>	<i>Terms</i>	<i>Media</i>	<i>Workflow</i>	<i>Algorithm</i>	<i>Combination</i>
Mean	500	20	20	50,000	3,000
Median	10	5	2	5,000	50
Deviation	30	5	5	200	20
Distribution	90	40	15	20	120
Correlation	15	10	5	20	90
Probability	140	15	20	50	150
Phonetics	50	5	10	20	50
Regular expr.	920	100	50	40	1,500
References	720	120	30	5	900
Association	610	60	10	5	420
UDC	530	120	20	5	660
Keywords	820	100	10	5	600
Translations	245	20	5	5	650
Corrections	60	10	5	5	150
External res.	40	30	5	5	40

Statistics methods have shown to be an important means for successfully optimising result matrices. The most widely implemented methods for the creation of result matrices are intermediate result matrices based on regular expressions and

intermediate result matrices based on combined regular expressions, classification, and statistics, giving their numbers special weight. Based on these per-instance numbers this results in demanding requirements for complex applications – On numerical data: Millions of calls are done per algorithm and dataset, hundreds in parallel/compact numeric routines. On “terms”: Hundred thousands of calls are done per sub-workflow, thousands in parallel/complex routines, are done.

Most resources are used for one application scenario only. Only 5–10 percent overlap between disciplines – due to mostly isolated use. Large benefits result from multi-disciplinary multi-lingual integration. The multi-lingual application adds an additional dimension to the knowledge matrix, which can be used by most discovery processes. As this implemented dimension is of very high quality the matrix space can benefit vastly from content and references.

VII. CONCLUSION AND FUTURE WORK

A number of structuring elements and workflow procedures have been successfully implemented for processing objects from knowledge resources, which allow optimising result matrices in very flexible ways.

First, long-term multi-disciplinary and multi-lingual knowledge resources can provide a solid source of structured content and references for a wealth of result matrices. The long-term results confirm that for the usability the organisation of the content and the data structures are most important and should have the overall focus compared to algorithm adaptation and optimisation. Nevertheless, the computational requirements may be very high but compared against the long-term data creation issues, they should be regarded secondary from the scientific point of view.

Second, employing a classification like UDC has shown to be a universal and most flexible solution with statistics for supporting long-term multi-disciplinary knowledge resources.

Third, computing optimised result matrices from objects of universally classified knowledge resources can be efficiently supported by various statistics and probability measures. With the quality and quantity of matrix elements this can also improve the decision making processes within the workflows.

The research conducted provided that advanced discovery will have to go into depth as well as into broad surface of the context of the multi-disciplinary and multi-lingual information in order to effectively improve the quality for most workflows. Many of these workflow processes can be very well parallelised on HEC resources. A typical case where parallelisation is required is the implementation of an application creating result matrices and used with many parallel instances. This introduces benefits for the applicability of the discovery facing big data resources to be included. The integration of the above strategies and means has proven an excellent method for computing optimised result matrices. Future work will be focussed on the workflow processes and standardisation and best practice for container and resources’ objects.

ACKNOWLEDGEMENTS

We are grateful to all national and international partners in the GEXI cooperations for their support and contributions.

Special thanks go to the scientific colleagues at the Gottfried Wilhelm Leibniz Bibliothek (GWLB) Hannover, especially Dr. Friedrich Hülsmann, for prolific discussion, inspiration, and practical case studies.

REFERENCES

- [1] A. Woodie, “Forget the Algorithms and Start Cleaning Your Data,” *Datanami*, 2014, March 26, 2014, URL: http://www.datanami.com/datanami/2014-03-26/forget_the_algorithms_and_start_cleaning_your_data.html [accessed: 2014-04-03].
- [2] C.-P. Rückemann, “Sustainable Knowledge Resources Supporting Scientific Supercomputing for Archaeological and Geoscientific Information Systems,” in *Proc. Third INFOCOMP 2013*, Nov. 17–22, 2013, Lisbon, Portugal, 2013, pp. 55–60, ISSN: 2308-3484, ISBN: 978-1-61208-310-0.
- [3] C.-P. Rückemann, “High End Computing for Diffraction Amplitudes,” in *Proceedings ICNAAM 2013*, Rhodes, Greece, vol. 1558. AIP Press, 2013, pp. 305–308, ISBN: 978-0-7354-1184-5, ISSN: 0094-243X, DOI: 10.1063/1.4825483.
- [4] U. Inden, D. T. Meridou, M.-E. C. Papadopoulou, A.-C. G. Anadiotis, and C.-P. Rückemann, “Complex Landscapes of Risk in Operations Systems Aspects of Processing and Modelling,” in *Proc. Third INFOCOMP 2013*, Nov. 17–22, 2013, Lisbon, Portugal, 2013, pp. 99–104, ISSN: 2308-3484, ISBN: 978-1-61208-310-0.
- [5] P. Leitão, U. Inden, and C.-P. Rückemann, “Parallelising Multi-agent Systems for High Performance Computing,” in *Proc. Third INFOCOMP 2013*, Nov. 17–22, 2013, Lisbon, Portugal, 2013, pp. 1–6, ISSN: 2308-3484, ISBN: 978-1-61208-310-0.
- [6] “LX-Project,” 2014, URL: <http://www.user.uni-hannover.de/cpr/xrprojs/en/#LX> (Information) [accessed: 2014-01-12].
- [7] 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.
- [8] “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/udccsummary/php/index.php> [accessed: 2014-01-12].
- [9] “UDC Online,” 2014, <http://www.udc-hub.com/> [acc.: 2014-01-12].
- [10] C.-P. Rückemann, “Archaeological and Geoscientific Objects used with Integrated Systems and Scientific Supercomputing Resources,” *International Journal on Advances in Systems and Measurements*, vol. 6, no. 1&2, 2013, pp. 200–213, ISSN: 1942-261x.
- [11] Y. Dodge, *The Oxford Dictionary of Statistical Terms*. Oxford University Press, 2006, ISBN: 0-19-920613-9.
- [12] B. S. Everitt, *The Cambridge Dictionary of Statistics*, 3rd ed. Cambridge University Press, Cambridge, 2006, ISBN: 0-521-69027-7.
- [13] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006, ISBN: 0-387-31073-8.
- [14] R. D. Drennan, *Statistics in Archaeology*, 2008, in: Pearsall, Deborah M. (ed.), *Encyclopedia of Archaeology*, pp. 2093-2100, Elsevier Inc., ISBN: 978-0-12-373962-9.
- [15] D. Lindley, “The Philosophy of Statistics,” *Journal of the Royal Statistical Society*, 2000, JSTOR 2681060, Series D 49 (3), pp. 293–337, DOI: 10.1111/1467-9884.00238.
- [16] “Universal Decimal Classification Consortium (UDCC),” 2014, URL: <http://www.udcc.org> [accessed: 2014-01-12].
- [17] “Universal Decimal Classification (UDC),” 2014, Wikipedia, URL: http://en.wikipedia.org/wiki/Universal_Decimal_Classification [accessed: 2014-01-12].
- [18] A. Slavic, “UDC libraries in the world - 2012 study,” *universaldecimalclassification.blogspot.de*, 2012, Monday, 20 August 2012, URL: <http://universaldecimalclassification.blogspot.de/2012/08/udc-libraries-in-world-2012-study.html> [accessed: 2014-01-12].
- [19] C.-P. Rückemann, “Integrating Information Systems and Scientific Computing,” *International Journal on Advances in Systems and Measurements*, vol. 5, no. 3&4, 2012, pp. 113–127, ISSN: 1942-261x.
- [20] “Creative Commons Attribution Share Alike 3.0 license,” 2012, URL: <http://creativecommons.org/licenses/by-sa/3.0/> [accessed: 2014-01-12].