

Krassimira Ivanova, Milena Dobрева, Peter Stanchev, George Totkov
(editors)

on and similar papers at core.ac.uk

provided by Bulgarian Digital Mat

Innovative Applications of Automated Metadata Generation

Plovdiv University Publishing House "Paisii Hilendarski"
2012, Plovdiv, Bulgaria

**Access to Digital Cultural Heritage:
Innovative Applications of Automated Metadata Generation**

Edited by:

Krassimira Ivanova, Milena Dobрева, Peter Stanchev, George Totkov

Authors (in order of appearance):

Krassimira Ivanova, Peter Stanchev, George Totkov, Kalina Sotirova, Juliana Peneva, Stanislav Ivanov, Rositza Doneva, Emil Hadjikolev, George Vragov, Elena Somova, Evgenia Velikova, Iliya Mitov, Koen Vanhoof, Benoit Depaire, Dimitar Blagoev

Reviewer: Prof., Dr. Avram Eskenazi

Published by: Plovdiv University Publishing House "Paisii Hilendarski"

2012, Plovdiv, Bulgaria

First Edition

The main purpose of this book is to provide an overview of the current trends in the field of digitization of cultural heritage as well as to present recent research done within the framework of the project D002-308 funded by Bulgarian National Science Fund. The main contributions of the work presented are in organizing digital content, metadata generation, and methods for enhancing resource discovery.

Printed in Bulgaria by Plovdiv University

24, Tsar Assen, Str., Plovdiv-4000, Bulgaria

All Rights Reserved

© This compilation: K. Ivanova, M. Dobрева, P. Stanchev, G. Totkov 2012

© The chapters: the contributors 2012

© The cover: K. Sotirova 2012

ISBN: 978-954-423-722-6

Plovdiv, 2012

Acknowledgements

This book summarises the outcomes of several recent research projects. It is quite a complex project in terms of scope and number of contributors. We would like to particularly thank all authors for their enthusiasm and commitment.

We would also like to thank all our colleagues who were supportive of our work and thus helped to develop our competences on the topic.

The projects which helped to develop our ideas and test them in real life were supported by:

- the Bulgarian National Science Fund, namely the Project D002-308 "MetaSpeed: Automated Metadata Generating for e-Documents Specifications and Standards";
- the Framework Programme 7 (FP7) of the European Commission, namely the project RI-246686 "OpenAIRE: Open Access Infrastructure for Research in Europe";
- the Hasselt University in Belgium, namely the Projects R-1875 "Search in Art Image Collections Based on Colour Semantics" and R-1876 "Intelligent systems' memory structuring using multidimensional numbered information spaces".

The involvement of Bulgarian researchers in these projects in particular was made possible through the cooperation of several institutions: we would like to express our gratitude to the *Institute of Mathematics and Informatics – Bulgarian Academy of Sciences, Plovdiv University "Paisii Hilendarski", New Bulgarian University, and Hasselt University in Belgium* for providing excellent conditions for collaboration.

A number of international as well as national events gave us possibilities to present our visions and discuss with other colleagues from different professional communities; these discussions were also of great contribution to our work. Of particular help were the international events organised by the Member of the European Parliament Emil Stoyanov; he also welcomed the idea of creating this collection.

We also want to thank our reviewer Professor Avram Eskenazi for his helpful remarks during the process of preparation of the content of this publication.

Last but not least, we would like to thank Emilia Todorova from Glasgow Caledonian University for the help with language revision and to Viktoria Naoumova from the Institute of Mathematics and Informatics for technical assistance.



Table of Contents

Acknowledgements.....	3
Table of Contents	5
List of Abbreviations	9
Introduction	13
Chapter 1: Digitization of Cultural Heritage – Standards, Institutions, Initiatives	23
1 Cultural Heritage.....	23
2 The Three Building Blocks of Digital Heritage.....	26
2.1 Digitization	26
2.2 Access	28
2.3 Preservation	29
3 The Importance of Metadata	31
4 Metadata Schemas and Standards Used in Cultural Heritage...33	
4.1 Common Standards.....	34
4.2 Standards for Resource Discovery	36
4.3 Specific Standards	37
4.4 Other Standards Relevant to Cultural Heritage	40
5 Digital Library	41
5.1 Basic Definitions	42
5.2 The Contemporary Models of Digital Libraries.....	43
5.3 Repository Software.....	50
6 Initiatives on World and European Level	52
6.1 Library and Scientific Open-access Initiatives	53
6.2 Examples of Initiatives that Change the Digital World.....	56
6.3 Initiatives, Connected with Data Content Standards	59
7 The User and the New Digital World	60
7.1 Users: between Policies and Real Involvement	61

7.2	User Involvement in Digital Libraries Development.....	62
7.3	User Studies	63
8	Conclusion	64
	Bibliography.....	65
Chapter 2: REGATTA – Regional Aggregator of Heterogeneous		
	Cultural Artefacts	69
1	Introduction	69
2	Aggregators of Digital Content for Cultural Artefacts in EU.....	71
3	The Prototype REGATTA–Plovdiv.....	72
3.1	The Functional Scheme of REGATTA	74
3.2	Data Model in REGATTA	75
3.3	Technological Aspects.....	78
4	Virtual Tours in REGATTA	80
4.1	Panoramic Virtual Tours.....	81
4.2	3D-Virtual Tours.....	82
5	Presentation of Plovdiv Ethnographic Museum in REGATTA	83
5.1	Movable Artefacts	83
5.2	Virtual Tours of the Plovdiv Ethnographic Museum	87
6	The Next Step – Enforcing the Data Management with Data Mining Tools	92
7	Conclusion	94
	Bibliography.....	94
Chapter 3: Automated Metadata Extraction from Art Images		
1	Introduction	97
2	Semantic Web	99
3	The Process of Image Retrieval	101
3.1	Text-Based Retrieval	101
3.2	Content-Based Image Retrieval (CBIR)	104
4	The Gaps	106
4.1	Sensory Gap	107
4.2	Semantic Gap	108
4.3	Abstraction Gap.....	109
4.4	Subjective Gap.....	110
5	User Interaction	111
5.1	Complexity of the Queries	111
5.2	Relevance Feedback.....	112
5.3	Multimodal Fusion.....	113
6	Feature Design	114

6.1	Taxonomy of Art Image Content	115
6.2	Visual Features.....	117
6.3	MPEG-7 Standard	123
7	Data Reduction	127
7.1	Dimensionality Reduction	127
7.2	Numerosity Reduction	134
8	Indexing	137
9	Retrieval Process.....	140
9.1	Similarity.....	140
9.2	Techniques for Improving Image Retrieval.....	146
10	Conclusion	147
	Bibliography.....	148

Chapter 4: APICAS – Content-Based Image Retrieval in Art

	Image Collections Utilizing Colour Semantics.....	153
1	Colour – Physiology and Psychology	153
1.1	Physiological Ground of the Colour Perceiving.....	155
1.2	Image Harmonies and Contrasts	157
1.3	Psychological Colour Aspects	159
2	Art Image Analyzing Systems	160
3	Proposed Features.....	163
3.1	Colour Distribution Features	164
3.2	Harmonies/Contrasts Features	166
3.3	Formal Description of Harmonies/Contrasts Features Using HSL- artist Colour Model.....	170
3.4	Local Features, based on Vector Quantization of MPEG-7 Descriptors over Tiles	176
3.5	Other Attributes	178
4	APICAS: The System Description	179
4.1	Functional Requirements.....	180
4.2	APICAS Architecture.....	181
4.3	APICAS Ground	183
4.4	APICAS Functionality	183
5	Experiments	192
5.1	Analysis of the Visual Features.....	192
5.2	Analysis of the Harmonies/Contrast Descriptors.....	194
5.3	Analysis of the Local Features.....	197
6	Conclusion	200
	Bibliography.....	201

Chapter 5: Automatic Metadata Generation and Digital Cultural Heritage	203
1 Automatic Generation of Metadata	203
1.1 Regular Expressions	204
1.2 Rule-based Parsers	204
1.3 Machine Learning Algorithms	205
2 Data Mining	205
3 Data Extraction from Web Documents Using Regular Expressions	209
3.1 Data Extraction by Learning Restricted Finite State Automata	210
3.2 Program Realization	213
3.3 Experiments	214
4 ArmSquare: an Association Rule Miner Based on Multidimensional Numbered Information Spaces	218
4.1 A Brief Overview of Previous ARM Algorithms	219
4.2 Association Rule Miner ArmSquare	221
4.3 Multidimensional Numbered Information Spaces	222
4.4 Algorithm Description of ArmSquare	223
4.5 Program Realization	227
4.6 Advanced Specifics of ArmSquare	229
4.7 Implementation	229
5 PGN: Classification with High Confidence Rules	232
5.1 The Structure of CAR-algorithms	233
5.2 Algorithm Description of PGN Classifier	235
5.3 PGN and Predictive Analysis in Art Collections	241
6 Metric Categorization Relations Based on Support System Analysis	246
6.1 The Semantic Complexity	246
6.2 Meta-PGN: Algorithm Description	247
6.3 Program Realization	248
6.4 The Next Step: Application in the Field	249
7 Conclusion	249
Bibliography	251

List of Abbreviations

5M	Multicultural, Multilingual, Multimodal, Multivariate, Modelling
5S	Streams, Structures, Spaces, Scenarios, and Societies
AAT	Art and Architecture Thesaurus
ACRI	Associative Classifier with Reoccurring Items
ACTA	Anti-Counterfeiting Trade Agreement
AIP	Archival Information Package
APICAS	Art Painting Image Colour Aesthetics and Semantics
ARC-AC	Association Rule-based Categorizer for All Categories
ARC-BC	Association Rule-based Categorizer By Category
ARM	Association Rule Mining
ArM	Archive Manager
ARUBAS	Association RULE BAsed Similarity framework
BIDL	Bulgarian Iconographical Digital Library
CAD	Computer-aided Design
CAR	Class-Association Rules
CATCH	Continuous Access to Cultural Heritage
CBA	Classification Based on Associations
CBIR	Content-Based Image Retrieval
CCA	Curvilinear Component Analysis
CCSDS	Consultative Committee for Space Data Systems
CDWA	Categories for the Description of Works of Art
CH	Cultural Heritage
CHO	Cultural and Historical Objects
CIDOC CRM	International Committee for Documentation – Conceptual Reference Model
CL	Colour Layout
CMAR	Classification based on Multiple Association Rules
CMY	Cyan-Magenta-Yellow
CONA	Cultural Objects Name Authority
CorClass	Correlated Association Rule Mining for Classification
CORDIS	Community Research & Development Information Service
CPAR	Classification based on Predictive Association Rules
CS	Colour Structure

CSDGM	Content Standard for Digital Geospatial Metadata
DACS	Describing Archives: a Content Standard
DC	Dominant Colour
DC	Dublin Core
DCP	Data Coverage Pruning
DELOS	Network of Excellence on Digital Libraries
DHO	Digital Humanities Observatory
DIP	Dissemination Information Package
DL	Digital Library
DLRM	Digital Libraries Reference Model
DOI	Digital Object Identifier
DWT	Discrete Wavelet Transform
EAD	Encoded Archival Description
EC	European Commission
ECDL	European Conference on Digital Libraries
EDL	European Digital Library
EDM	Europeana Data Model
EH	Edge Histogram
EMD	Earth Mover's Distance
EOF	Empirical Orthogonal Function
Fedora	Flexible Extensible Digital Object Repository Architecture
FOIL	First Order Inductive Learner
FP7	Seventh Framework Programme
FRBR	Functional Requirements for Bibliographic Records
FRBROO	FRBR – Object Oriented
GIS	Geographic Information System
GLAM	Galleries, Libraries, Archives, Museums
GLOH	Gradient Location and Orientation Histogram
GPS	Global Positioning System
HARMONY	Highest confidence cAssification Rule Mining fOr iNstance-centric classifYing
HSIS	Humanities Serving Irish Society
HSL	Hue-Saturation-Luminance
HSV	Hue-Saturation-Value
HT	Homogeneous Texture
HTML	Hyper-Text Markup Language
ICCROM	International Centre for the Study of the Preservation and Restoration of Cultural Property
ICOM	International Council of Museums
ICT-CIP	Information and Communication Technologies – Competitiveness and Innovation Framework Programme
IDABC	Interoperable Delivery of European eGovernment Services to public Administrations, Businesses and Citizens
IFLA	International Federation of Library Associations

IMI-BAS	Institute of Mathematics and Informatics – Bulgarian Academy of Sciences
IRI	International Resource Identifier
ISAAR(CPF)	International Standard Archival Authority Record for Corporate Bodies, Persons and Families
ISAD(G)	General International Standard Archival Description
ISOC	Internet Society
IT	Information Technology
JISC	Joint Information Systems Committee
LESH	Local Energy based Shape Histogram
LIDAR	Light Detection And Ranging
LIDO	Light Information Describing Objects
LLE	Locally Linear Embedding
MARC	MAchine-Readable Cataloging
MBR	Minimum Bounding Rectangle
MDS	Multi Dimensional Scaling
MET	Metropoliten Museum of Art
METS	Metadata Encoding & Transmission Standard
MINERVA	MInisterial NEtwoRk for Valorising Activities in digitisation
MODS	Metadata Object Description Schema
MPEG	Moving Picture Experts Group
NSDL	National Science Digital Library
NURBS	Non Uniform Rational BSpline
OAI-PMH	Open Archives Initiative Protocol for Metadata Harvesting
OAIS	Open Archival Information System
OpenAIRE	Open Access Infrastructure Research for Europe
ORE	Ontology Rule Editor
OWL	Web Ontology Language
PCA	Principal Component Analysis
PGN	Pyramidal Growing Network
PP	Projection Pursuit
R&D	Research and Development
RDF	Resource Description Framework
RDFS	Resource Description Framework Schema
REGATTA	REGional Aggregator of heTerogeneous culTural Artefacts
RGB	Red-Green-Blue
RYB	Red-Yellow-Blue
SAIL	Semi-Automated Interactive Learning systems
SC	Scalable Colour
SDA	Symbolic Data Analysis
SGML	Standard Generalized Markup Language
SIFT	Scale-Invariant Feature Transform
SIP	Submission Information Package
SRES	Self-supervised web relation Extraction System
SRSWOR	Simple Random Sample WithOut Replacement

SRSWR	Simple Random Sample With Replacement
SURF	Speeded Up Robust Feature
SVD	Singular Value Decomposition
SVM	Support Vector Machines
TEL	The European Library
TFPC	Total From Partial Classification
TGN	Thesaurus of Geographic Names
TPDL	Theory and Practice of Digital Libraries
ULAN	Union List of Artist Names
UNESCO	United Nations Educational, Scientific and Cultural Organization
URI	Uniform Resource Identifier
URL	Uniform Resource Locator
URN	Uniform Resource Name
VQ	Vector Quantization
VRA	Visual Resources Association
W3C	World Wide Web Consortium
WDL	World Digital Library
WIPO	World Intellectual Property Organization
XML	eXtensible Markup Language