fi yuo cna raed tihs, yuo hvae a sgtrane mnid. Cna yuo raed tihs? Olny 55 plepoe out of 100 can

# Making your Metadata Beautiful

13 July 2016

Heather Rosie & Sara Gould

The British Library e-theses online service (EThOS)

# What is 'beautiful' metadata?

## "ACRC" ...

- **Accurate**: correct transcription; no duplication; no typos

- **Consistent**: same type of data recorded in the same place; data presented in the same way

- **Rich**: as much relevant data as possible, e.g. abstracts, identifiers, funder, supervisor

- **Current**: material available now needs to be discoverable now

# Benefits of beautiful ("ACRC") metadata

- Optimises resource discovery – and increases usage

- Maximises workflow efficiencies

- Supports accurate and reliable statistics/metrics (e.g. number of theses awarded in the UK by institution, subject or date)

- Increases interoperability - allowing data from multiple sources to be combined

# Challenge 1- record sources

1. British Library catalogue records

2. University catalogue records

3. Institutional repository harvested records

4. BL Cataloguer created records

5. BL Imaging Team created records

# Challenge 2 - metadata formats

1. MARC Exchange Format (ISO 2709)

2. MARC XML

3. OAI_DC XML

4. UKETD_DC XML

5. RIOXX

6. Tab-delimited (Excel)

# Challenge 3 - skills/knowledge of the people creating metadata

1. Cataloguers

2. Students

3. Repository staff (managers, admin, clerical)

# Challenge 4 - processing and configuration of data storage and export

1. Programmers

2. Repository Managers

3. Administrative staff

4. EThOS Metadata Manager

# Challenge 5 - special characters/diacritics

1. Spelt out ('alpha')

2. UTF-8 character set

3. Local convention (e.g. ^ for superscript number)

4. Data represented in $ strings

5. Characters omitted, garbled or replaced by ??

# All these challenges together can produce 'poor' metadata

- Character encoding (MARC-8 v. UTF-8)

- Typographical errors

- Scanned texts (machine error)

- Missing data (abstracts, subtitles, authors, language, subjects, qualification name)

- No standardisation (e.g. qualification name; author names; representation of words/symbols)

- Lack of identifiers

- Duplication

# Character encoding

- Web standard is UTF-8 (Unicode);   MARC character set = MARC-8

- Software applications may not recognise 'foreign' character sets and will display as strange characters and/or question mark(s)

- Some characters are 'invalid' in XML (and some invalid characters are invisible, e.g. 'null'); web browsers will not display XML documents containing invalid characters

- Use of invalid characters impacts on resource discovery and display and, in some cases, can cause software failure

# Diacritics in title

The University of Wales, Lampeter

A study and edition of Ima&#772;m Abd al-Azi&#772;z
b. Ali&#772; b. al-Izz al-Baghda&#772;di&#772;
al-Bakri&#772; al-H&#803;anbali&#772;
al-Maqdisi&#772; Junnat
al-S&#803;a&#772;biri&#772;n al-Abra&#772;r Wa
Jannat al-Mutawakkili&#772;n al-Akhya&#772;r

Author:                  Al-Olabi, Adnan al-Hamwi
Awarding Institution:    The University of Wales, Lampeter
Current Institution:     The University of Wales Trinity Saint David
Awarded:                 2003

Thesis available for immediate download

Advisor:              Not available          Sponsor:
Qualification name:   PhD                    Qualification Level:
EThOS Persistent ID:  uk.bl.ethos.503578

# Amended diacritics in title

| Title: | A study and edition of Imām Abd al-Azīz b. Alī b. al-Izz al-Baghdādī al-Bakrī al-Ḥanbalī al-Maqdisī Junnat al-Ṣābirīn al-Abrār Wa Jannat al-Mutawakkilīn al-Akhyār | | |
|---|---|---|---|
| Author: | Al-Olabi, Adnan al-Hamwi | | |
| Awarding Body: | The University of Wales, Lampeter | | |
| Current Institution: | University of Wales Trinity Saint David | | |
| Date of Award: | 2003 | | |
| Availability of Full Text: | Access through EThOS: | Thesis available for immediate download. Please login/register to view download & delivery options. | |
| EThOS Persistent ID: | uk.bl.ethos.503578 | | |
| Supervisor: | Not available | Sponsor: | JISC Digital Islam |
| Qualification Name: | Thesis (Ph.D.) | Qualification Level: | Doctoral |
| Abstract: | | | |
| No abstract available | | | |
| Share: | ShareThis   Facebook   Tweet   LinkedIn   Email   CiteULike   Blogger | | |

# Typographical errors (human)

Causes include:

- Focusing on quantity rather than quality due to work pressures/targets

- Lack of interest/boredom

- Seeing what you think you see (SWYTYS) rather than what is actually there

fi yuo cna raed tihs, yuo hvae a sgtrane mnid too. Cna yuo raed tihs? Olny 55 plepoe out of 100 can

i cdnuolt blveiee taht I cluod aulaclty uesdnatnrd waht I was rdanieg. The phaonmneal pweor of the hmuan mnid, aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it dseno't mtaetr in waht oerdr the lttere s in a wrod are, the olny iproamtnt tihng is taht the frsit and lsat ltteer be in the rghit pclae.

The rset can be a taotl mses and you can sitll raed it whotuit a pboerlm. Tihs is bcuseae the huamn mnid deos not raed ervey lteter by istlef, but the wrod as a wlohe. Azanmig huh? Yaeh, and I awlyas tghuhot slpeling was ipmorantt! If you can raed this bceome a fna.

# The 'SWYTS' effect

| | |
|---|---|
| special reference | specific reference |
| private | pirate |
| conservation | conversation |
| hypertensive | hypersensitive |
| isoform | inform |
| economic aspects | economic analysis |
| biochemical | biomechanical |
| experiential | experimental |

# Title errors

Awarenees in ageing

Awareness in ageing

Efficient parallel genetic algorithms applied to numerical optimisation

Efficient Parallel Genetic Alogrithms Applied to Numerical Optimisation

The doll : the figure of the doll in culture and theory

The figure of the doll in culture and theory

# OCR errors (machine)

- Scanning texts is imperfect

- Special characters are often not recognised

- Words/sentences may be concatenated

- Review by humans is subject to the "SWYTYS" effect

# OCR errors in abstract

±a|$' dissertation is concerned primarily with interpreting     iD.H. Lawrence's poems,
rather than with imposing onto his poetry a     r pre-conceived critical hypothesis. My
method has been to start from the poetic texts and work outwards, producing a deliberately
intensive reading of the poems, centred on Lawrence, and not a wide-rangingcomparative study.
rArising tn~t of this close examination of the verse itself, is a         r discussion of the
relationship between Lawrence's art and hisihought,which I consider further in a short
introduction. It is a subject    r which I do not think has been discussed in relation to
Lawrence's poetry before. Briefly, I find that it is the art, and specifically the poetry,
which has primacy, and see Lawrence's development as one in which what are initially poetic
formulations and artistic choices, are elevated to the level of conscious philosophic belief.
I suggest how, in Lawrence's bad, sermonizing verse, this process of composition is reversed,
and ideas forced back onto a poetrj from which they derive.In ay first chapter, I deal with
Lawrence's early verse in sucha way as to suggest how, in its materials and methods, later
developments    rmay be seen in embryo. Then in four chapters dealing with all Lawrence' a
r major verse, I follow the different relationship between thought and poetry at different
stages of Lawrence's career. My effort in each chapter is to bring out the distinctive
character of each book of poems discussed, while having an eye also for the sort of
continuity I see in the whole. At the end of W,dissertation, I have included what is intended
to be a complete catalogue of the criticism of Lawrence's poetry in English, arranged
chronologically.

# Special characters (e.g. beta)

Many ways to write special characters:

- Spelt out (beta)

- Roman alphabet letter rather than Greek alphabet (b)

- Image rather than text (β)

- Character entity reference (&#946;)

- UTF-8 value (β)

# Encodings for Greek beta (fileformat.info)

| | |
|---|---|
| HTML Entity (decimal) | &#946; |
| HTML Entity (hex) | &#x3b2; |
| HTML Entity (named) | &beta; |
| How to type in Microsoft Windows | Alt **+3B2** |
| UTF-8 (hex) | 0xCE 0xB2 (ceb2) |
| UTF-8 (binary) | 11001110:10110010 |
| UTF-16 (hex) | 0x03B2 (03b2) |
| UTF-16 (decimal) | 946 |
| UTF-32 (hex) | 0x000003B2 (3b2) |
| UTF-32 (decimal) | 946 |
| C/C++/Java source code | "\u03B2" |
| Python source code | u"\u03B2" |
| | More... |

# Example of special character errors

**Title loaded to EThOS (from Aleph)**:

Growth and characterisation of terrace graded virtual substrates with
SiÃ¢â‚¬Å¡ÂÃ¢â‚¬Å¡â€¹Ã‹Â£GeÃ‹Â£ 0.15 Ã¢â€°Â¤ x Ã¢â€°Â¤ 1

**Harvested title**:

Growth and characterisation of terrace graded virtual substrates with Si[subscript 1-x]Ge[subscript x] 0.15 ? x ? 1

**Correct title (i.e. utf-8 values)**:

Growth and characterisation of terrace graded virtual substrates with $Si_{1-x}Ge_x$ $0.15 \leq x \leq 1$

# How to create beautiful (ACRC) metadata

1. Quality control is essential

2. Combination of human and machine review

3. Software helps with the task

e.g. MARC Report

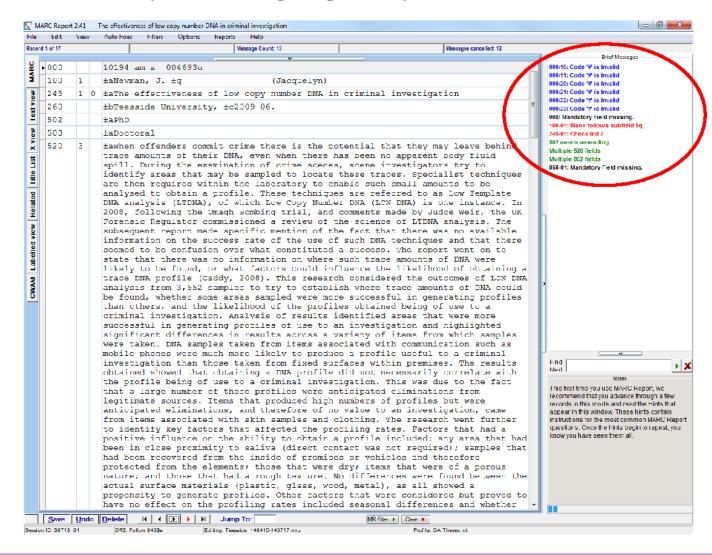http://www.marcofquality.com/soft/softindex.html

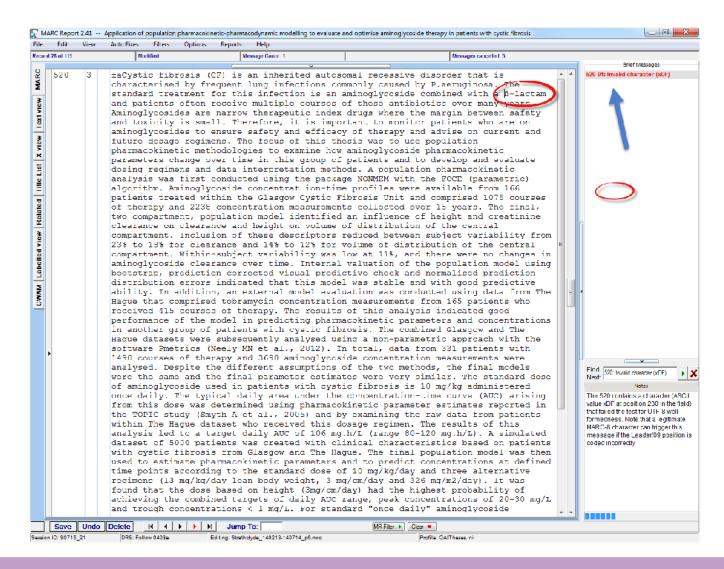# MARC Report/MARC Global

MARC Report can be used for:

- Checking errors (including invalid utf-8 characters)

- Creating a 'match key' and title list for de-duplication

- Converting data from XML to MARC and vice versa
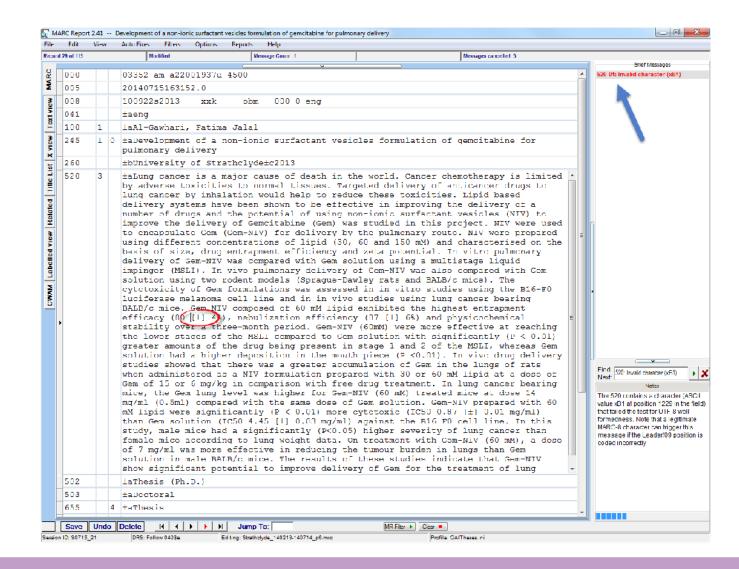
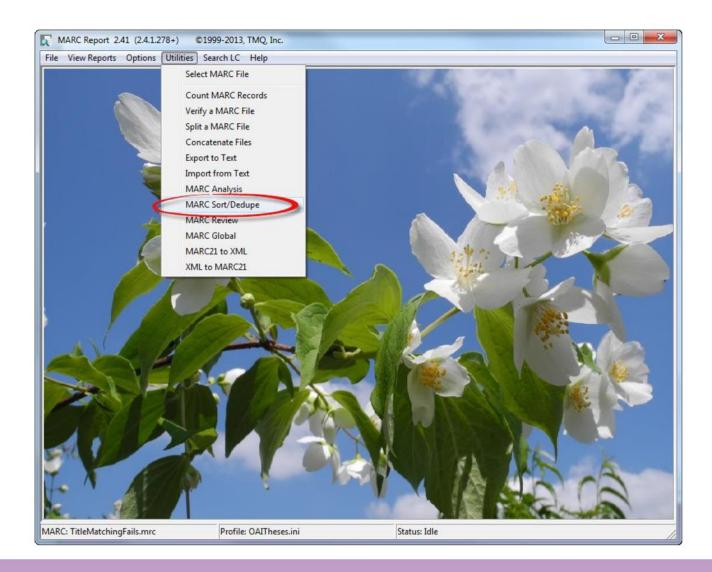- Lots of other applications!

# Data correction (errors highlighted)

# Unicode errors: Invalid Greek letter

# Unicode errors: Invalid maths symbol



www.bl.uk

28

# Deduplication

# Deduplication

- Using MARC Report Sort/deduplication - create a 'match key' that excludes punctuation, spacing, and diacritics.

- Author match on surname only (due to differences in recording forenames) and sort by title match key

# Author match title list

| RSN | 852 j | 245 ahnpb | 100 a |
|---|---|---|---|
| 1 | $j446687 | $aAn Investigation into the Longitudinal Rolling of Tubes Through Two grooved Rolls | $aAbdel-Haleem, A. M. S. |
| 2 | $j544618 | $aAn investigation into the longitudinal rolling of tubes through two grooved rools | $aAbdel-Haleem, Abdel R. M. S. |
| 3 | $j236625 | $aAn investigation of the emerging system of host governmental participation in the international oil industry with part | $aAbozrida, M. A. |
| 4 | $j544958 | $aAn investigation of the emerging system of host government participation in the international oil industry with partic | $aAbozrida, Mokhtar A. |
| 5 | $j294603 | $aThe polymerisation of lactic acid anhydrosulphite by anionic initiators | $aAdams, Luke Richard |
| 6 | $j544709 | $aThe polymerization of lactic acid anhydrosulphite by anionic initiators | $aAdams, L. R. |
| 7 | $j545069 | $aThe antigenic composition of Streptococcus faecalis associated with ineffective endocarditis | $aAitchison, Eileen J. |
| 8 | $j381964 | $aThe antigenic composition of Streptococcus faecalis associated with infective endocarditis | $aAitchison, E. J. |
| 9 | $j544866 | $aThe mechanisms of sulfur-containing metal complexes as UV-stabilisers | $aAl-Malaika, Sahar N. |
| 10 | $j344200 | $aThe mechanisms of sulphur-containing metal complexes as UV-stabilisers | $aAl-Malaika, S. N. |
| 11 | $j544883 | $aBound antioxidants in elastomers | $aAl-Mehdawe, Mohammed S. A. |
| 12 | $j237120 | $aBound antioxidants in electromers | $aAl-Mehdawe, M. S. A. |
| 13 | $j544770 | $aTheoretical studies of the mid-latitude ionosphere | $aAl-Naghmoosh, Ali A. |
| 14 | $j237163 | $aTheroetical studies of the mid-latitude ionosphere | $aAl-Naghmoosh, A. A. |
| 15 | $j312001 | $aStudies on the thermal decomposition behavior, kinetics and electrical conductivity of the non-isothermal decomp | $aAl-Sousi, Ghareeb Nemir |
| 16 | $j544689 | $aStudies on the thermal decomposition behaviour, kinetics and electrical conductivity of the non-isothermal decom | $aAl-Sousi, Ghareeb N. |
| 17 | $j307367 | $aBlock co-polymerization by transformation actions | $aAmass, Dorothy Gwendoline |
| 18 | $j544697 | $aBlock co-polymerization by transformation reactions | $aAmass, Dorothy G. |
| 19 | $j235168 | $aThe effect of growth conditions on #beta#-lactam resistance in Enterobacter cloacae | $aAnderson, E. M. |
| 20 | $j545046 | $aThe effect of growth conditions on B-lactam resistance in Enterobacter cloacae | $aAnderson, Elaine M. |
| 21 | $j544757 | $aThe effects of a dietary bacterial protein on mineral balance in rainbow trout (S. gairdneri Rich.) | $aAnglesea, Jonathan D. |
| 22 | $j331907 | $aThe effects of a dietary bacterial protein on mineral balance in rainbow trout (S. garidneri Rich.) | $aAnglesea, J. D. |
| 23 | $j237483 | $aA study of the scatetring of fast neutrons in large samples | $aAnvarian, S. P. T. |
| 24 | $j544965 | $aA study of the scattering of fast neutrons in large samples | $aAnvarian, Sattar P. T. |
| 25 | $j378554 | $aPhysicochemical characteristics of chlorofluorocarbon based inhalation aerosols | $aAshurst, I. C. |
| 26 | $j545065 | $aPhysicochemical characteristics of chlorofluorohydrocarbon based inhalation aerosols | $aAshurst, Ian C. |
| 27 | $j237582 | $aAtimulus-mitosis in the rat thymic lymphocyte | $aAtkinson, M. J. |
| 28 | $j544756 | $aStimulus-mitosis coupling in the rat thymic lymphocyte | $aAtkinson, Michael J. |
| 29 | $j331325 | $aDesign of corporate planning systems | $aBahrami, H. |
| 30 | $j544956 | $aDesign of corporate planning systems : development of a 'design framework' and its application to a specific sett | $aBahrami, Homa |
| 31 | $j253713 | $aThe effects of persistent anticholinesterase action at the neuromuscular junction | $aBamforth, John Philip |
| 32 | $j545074 | $aThe effects of persistent articholinesterase action at the neuromuscular junction | $aBamforth, John P. |
| 33 | $j448945 | $aNon-linear oil film force coefficients for a journal bearing operating under aligned and misaligned conditions | $aBannister, R. H. |

# Top 10 tips for metadata creation

1. Ensure accurate transcription

2. Use controlled vocabularies / standardised data entry

3. Use full names for authors, with correct entry element

4. Use the UTF-8 (Unicode) character set; avoid html markup, especially in the title:

   Preparation of <span style='font-family:Symbol'>h</span><sup>3</sup>-allymolybdenum complexes using <i>cis</i>-Mo(CO)<sub>4</sub></sub

5. Use sentence casing not ALL CAPS:

   CodeZebraOS (CODEZEBRAOS)

   GaAs (GAAS)

   MARS (Mars)

# Top 10 tips for metadata creation

6. Use correct punctuation and spacing:

    I like cooking my family and my pets

    I like cooking, my family, and my pets

7. Subject classify all theses

8. Always review OCR'd text

9. Use identifiers where possible

10. Use repeated fields for data of the same type (e.g. subject keywords, multiple authors)

# Test !

- A

- C

- R

- C

# Test !

- **Accurate**

- **C**

- **R**

- **C**

# Test !

- **Accurate**

- **Consistent**

- **R**

- **C**

# Test !

- **Accurate**

- **Consistent**

- **Rich**

- **C**

# Test !

- Accurate

- Consistent

- Rich

- Current

Heather.Rosie@bl.uk

EThOS.bl.uk