

NETSL  
Annual Spring Conference

April 15, 2010

# Mapping Bibliographic Metadata

**Carol Jean Godby, Research Scientist**



**OCLC™**

The world's libraries.  
Connected.

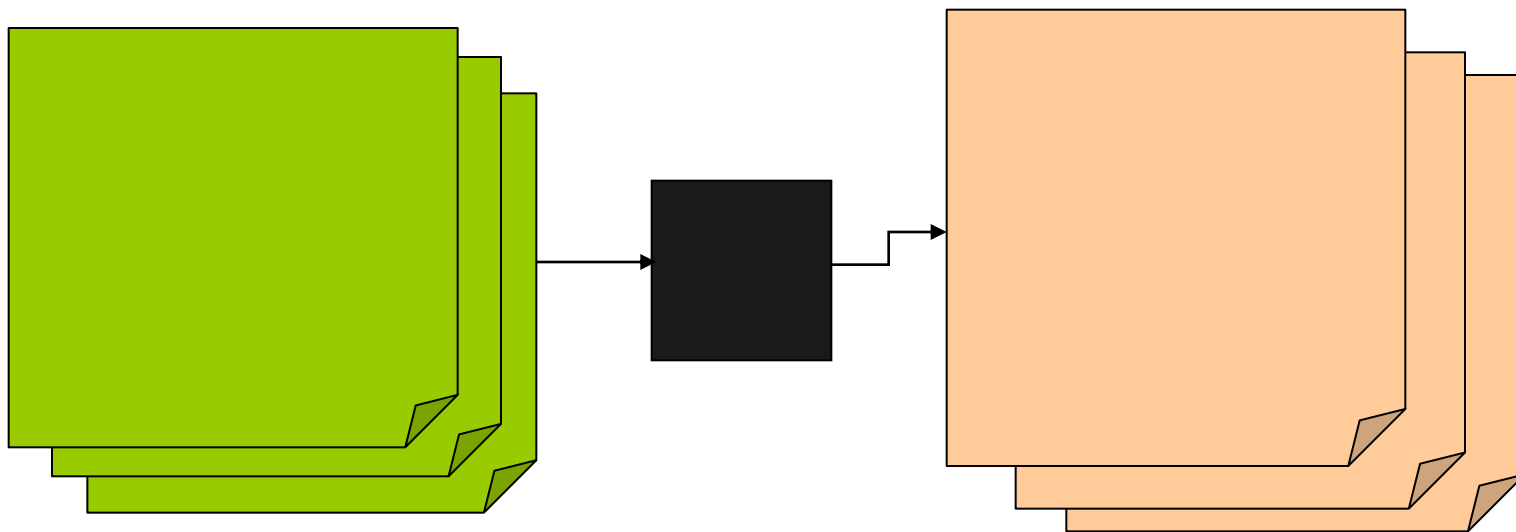
# Agenda for today



- 1. Crosswalks and how they are used at OCLC
- 2. A report of successes on some difficult problems
- 3. What is still unresolved
- 4. Some topics for group discussion

# The problem

My records are in this format...



...but they need to be in this other format.

# The supported translations

## Inputs

MARC 21-  
2709

ONIX Books

MARC XML

MODS

DC XML

OAI-DC XML

OCLC CDF

DC-Qualified

ONIX Serials

**OCLC's  
Common  
Data  
Format**

MARC 21-  
2709

## Outputs

OCLC MARC

OCLC CDF

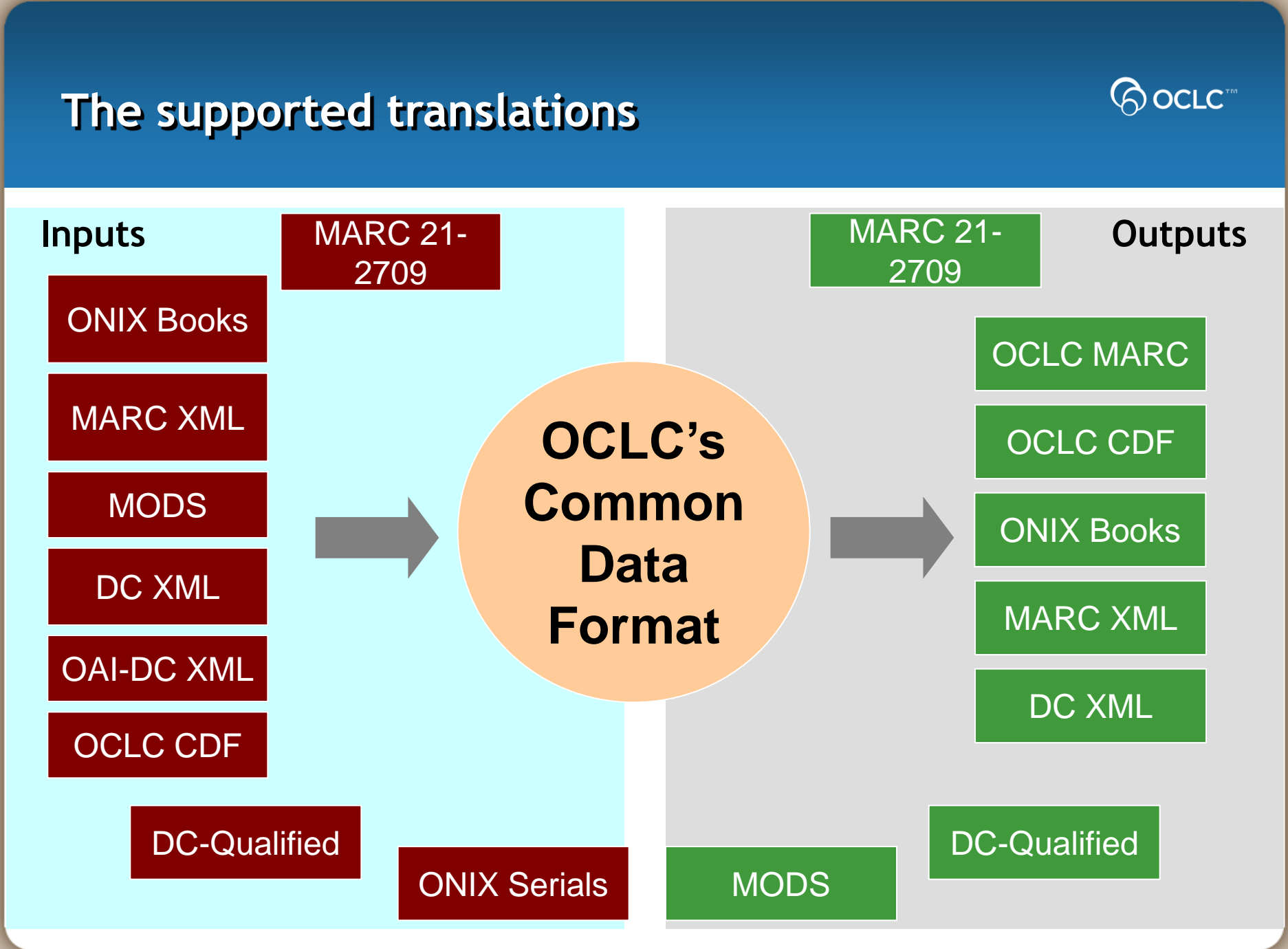
ONIX Books

MARC XML

DC XML

MODS

DC-Qualified



# What is a crosswalk?

“Crosswalks are used to ‘translate’ between different metadata element sets. The elements (or fields) in one metadata set are correlated with the elements of another metadata set that have the same or similar meanings. This is also sometimes called ‘semantic mapping.’”

# ONIX and MARC records



**<Product>**  
  **<RecordReference>**0892962844**.**  
  **<ProductIdentifier>**  
    **<ProductIDType>**03 **</>**  
    **<IDValue>**9780892962846**</>**  
  **</ProductIdentifier>**  
  **<ProductForm>**BB**</>**  
  **<Title>**  
    **<TitleType>**01 **</>**  
    **<TitleText>**McBain's Ladies**</>**  
  **</Title>**  
  **<Contributor>**  
    **<ContributorRole>**A01 **</>**  
    **<PersonNameInverted>**Hunter, Evan**</>**  
  **</Contributor>**  
  **<Subject>**  
    **<SubjectSchemeIdentifier>**02**</>**  
    **<SubjectHeadingText>**  
      Policewomen--Fiction.

Leader 00000 jm a22000005 4500  
008 g eng  
020 0892962844  
024 3# 9780892962846  
100 \$a Hunter, Evan  
245 \$a McBain's ladies  
260 \$b Mysterious Press \$d  
1988  
300 \$a 320 p.  
650 #2 \$a Policewomen -- Fiction

# ONIX is the international standard for the book industry



The screenshot shows a web browser window with the URL <http://www.editeur.org/12/About-Release-3.0/>. The website has a dark blue header with the 'EDITEUR' logo and a navigation menu. A sidebar on the left contains a 'STANDARDS' section with a dropdown for 'ONIX for Books' and a list of links including 'Overview', 'About Release 3.0', 'Release 3.0 Downloads', 'Code Lists', 'Previous Releases', and 'Agency terms in ONIX'. The main content area features a large banner for 'ONIX 2.0' with a background of binary code. Below the banner, the 'About Release 3.0' page is displayed, starting with the text: 'Release 3.0 is a major new version of the ONIX for Books standard, the first since 2001 that is not backwards-compatible with its predecessors. This extensive revision of the format has had two key drivers: the need to improve the handling of digital products (Guidance notes on describing digital products in ONIX are available), and the recognition that the price of maintaining backwards-compatibility has been the increasing number of 'deprecated' elements that have had to be maintained - and supported by ONIX receivers - even though they are no longer recommended for use.' The page also mentions improvements in other areas and provides links to download schemas, documentation, and sample messages.

EDITEUR

STANDARDS

ONIX

ONIX for Books

- Overview
- [About Release 3.0](#)
- Release 3.0 Downloads
- Code Lists
- Previous Releases
- Agency terms in ONIX

NEW

- Maintenance and support
- FAQs
- ONIX and MARC21

ONIX for Serials

Licensing Terms

FTP Filenaming

ONIX Identifier

Registration Formats

You Are Here: Home / Standards / ONIX / ONIX for Books / About Release 3.0

## About Release 3.0

Release 3.0 is a major new version of the ONIX for Books standard, the first since 2001 that is not backwards-compatible with its predecessors. This extensive revision of the format has had two key drivers: the need to improve the handling of digital products (Guidance notes on describing digital products in ONIX are available), and the recognition that the price of maintaining backwards-compatibility has been the increasing number of 'deprecated' elements that have had to be maintained - and supported by ONIX receivers - even though they are no longer recommended for use.

At the same time, the opportunity has been taken to introduce important improvements in other areas, although there are many data element groups where little or no change has been considered necessary.

For full background, an overview of the message structure, and a summary of key differences between Release 2.1 and Release 3.0, [read the Introduction to ONIX for Books 3.0](#).

To download the schemas, detailed documentation, and supplementary guidelines, go to the [Release 3.0 Downloads page](#).

To download a sample message as a zip file containing the message both as XML and as an annotated PDF,

# OCLC's ONIX to MARC Crosswalk



D2					fx																			
A B C D					E					F					G									
1					ONIX Reference Name					Notes/Comments					MANDATORY/OPTIONAL					MARC21 Data Element				
2																								
3					<Header>																			
4					<FromSAN>					SAN of the data sender. This must be the same of the sendingPartyID in the message data element in the transmission header.					Mandatory									
5					<SentDate>					Format is CCYYMMDDHHMM. System will generate.					Mandatory									
6					<MessageNumber>					Uniquely identifies the message. If not transmitted, the system will generate.					Mandatory					981 \$c				
7					<DefaultLanguageOfText>					List 74. Indicates the default language which is assumed for the text of products listed in the message, unless explicitly stated otherwise in a <Language> composite in the product record.					Optional					008/35-37				
8					<DefaultPriceTypeCode>					List 58. Indicates the default price type which is assumed for prices listed in the message, unless explicitly stated otherwise in a <Price> composite in the product record.					Optional									
9					<DefaultCurrencyCode>					List 96. Indicates the currency which is assumed for prices listed in the message, unless explicitly stated otherwise in a <Price> composite in the product record.					Optional									
10																								
11					<Product>										Mandatory									
12					<RecordReference>					The vendor control number (VCN) or unique identifier for the title.					Mandatory					001 037 \$a				
					<NotificationType>					List 1. 01 -- early notification 02 -- advance notification 03 -- notification confirmed from book-in-hand 04 -- update					Mandatory					<NotificationType> EncLvl 01 3 02 5 03 8 04 0				



## • OCLC Services

[Contract Cataloging  
Services for Publishers](#)[OCLC Metadata  
Services for Publishers](#)

- [Overview](#)
- [Ordering](#)
- [Support](#)

[See the OCLC press  
release.](#)

OCLC Services : OCLC Metadata Services for Publishers

## OCLC Metadata Services for Publishers

Enhanced ONIX metadata for publishers

### AT A GLANCE

- **Improves discoverability**—through the addition of subject analysis, classification numbers, reviews and more
- **Improves efficiency**—in creating and distributing metadata to supply chain partners
- **Increases marketability**—through the addition of subject analysis, classification numbers, evaluative content, reviews and more

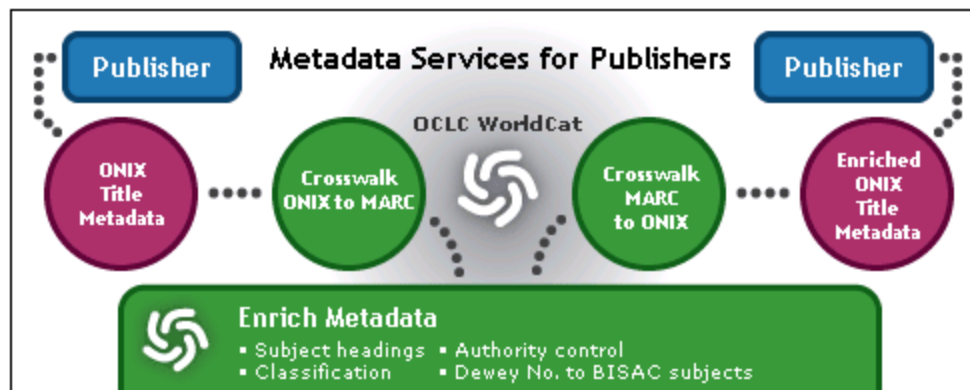
[View complete Overview >>](#)[Download the brochure >>](#)

### NEWS

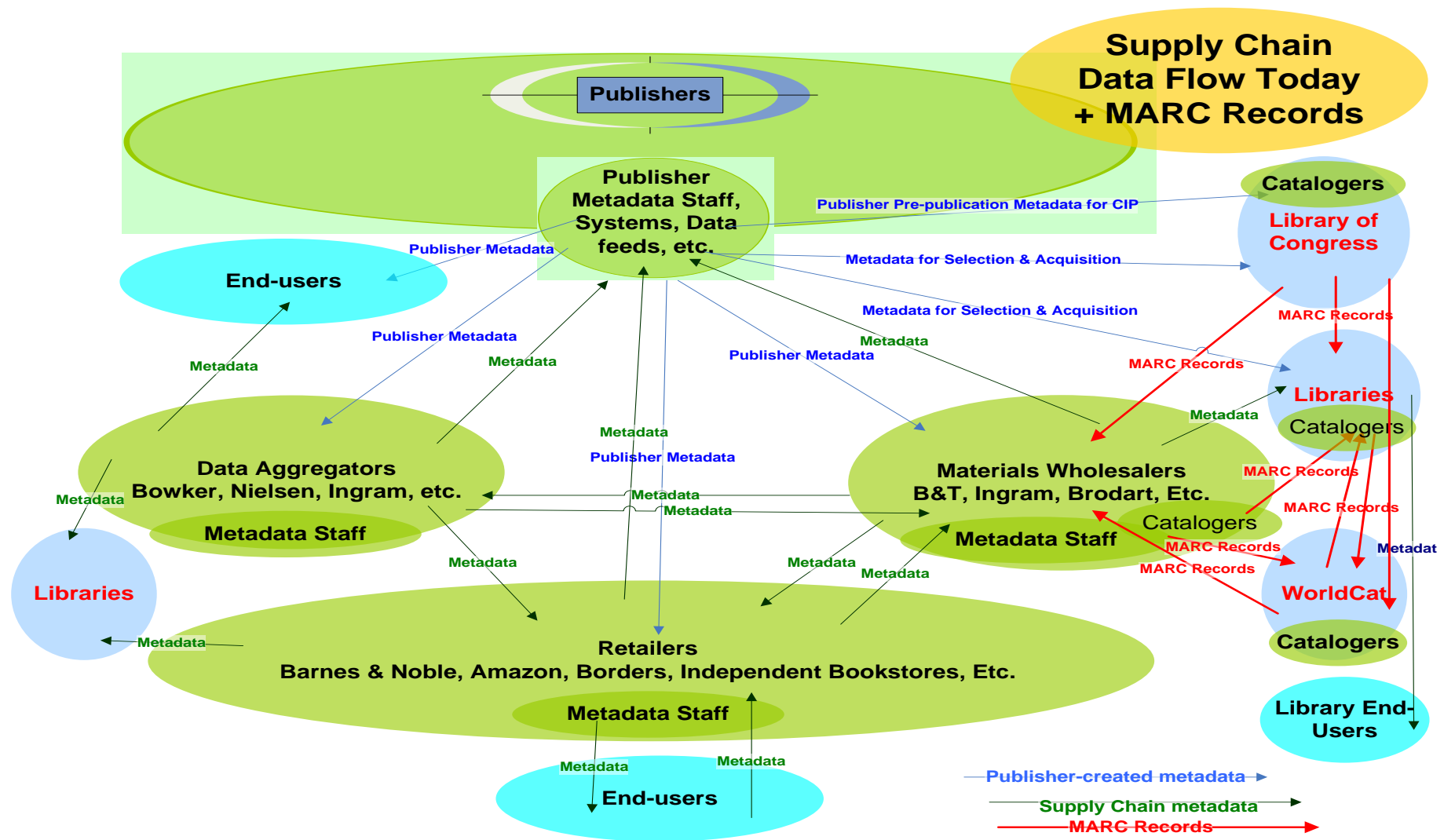
- [Webinar introduces Metadata Services for Publishers](#)

[More news >>](#)

OCLC's Metadata Services for Publishers can help you sell more titles by accepting your title metadata in ONIX format, enriching it by mining the WorldCat database, and delivering an enhanced ONIX file ready to use in supply chain systems and marketing communications.

[Learn how to order or request a quote for OCLC Metadata Services for Publishers >>](#)

## + MARC Records



# Metadata Services for Publishers: Highlights



- **ONIX records are obtained from publishers and mapped to MARC.**
- **If a match is found in WorldCat, fields from the corresponding MARC record are applied to it.**
- **If no match is found, the record is populated from fields in the closest FRBR cluster.**
- **The result is mapped back to ONIX to be delivered to publishers. It is also made available to the library community as an enhanced MARC record.**

OC LC 48893831 No holdings in OCL - 6 other holdings

Books Rec stat c Entered 20020130 Replaced 20070315095236.1

Type a ELvl l Srce d Audn Ctrl Lang eng  
BLvl m Form Conf 1 Biog MRec Ctry ja  
Cont b t GPub LitF 0 Indx 1  
Desc a Ills a Fest 0 DtSt s Dates 2001

040 OCC #c OCC #d OCL #d OCLCQ #d TEF #d OCL

020 4924600989

020 9784924600980

041 0 eng #a jpn

082 0 4 025.3/4 #2 22

090 Z699.35.M28 #b l58 2001

092 #b

049 OCLC

111 2 International Conference on I

245 1 0 Proceedings of the Internation  
Informatics ... [et al. ; edited by

246 1 #i At head of title: #a DC-2001

260 Tokyo : #b National Institute o

300 xi, 297 p. : #b ill. ; #c 30 cm.

546 English and Japanese.

<?xml version="1.0" encoding="UTF-8" ?>

- <dcterms>

- <qualifieddc xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:dcterms="http://purl.o

xmlns:dcmitype="http://purl.org/dc/dcmitype/" xmlns:xsi="http://www.w3.org/2001

xsi:noNamespaceSchemaLocation="http://dublincore.org/schemas/xmls/qdc/2006/0

<dc:contributor>Oyama, Keizo.</dc:contributor>

<dc:contributor>Gotoda, Hironobu.</dc:contributor>

<dc:contributor>Kokuritsu Jōhōgaku Kenkyūjo.</dc:contributor>

<dc:contributor>Dublin Core Metadata Initiative.</dc:contributor>

<dc:creator>International Conference on Dublin Core and Metadata Applications (200

<dcterms:issued>[2001]</dcterms:issued>

<dc:description>Includes bibliographical references.</dc:description>

<dc:description>Material support for the DC-2001 conference provided by: National In  
Dublin Core Metadata Initiative (DCMI), Japan Science and Technology Corporatio  
Information Science (ULIS), Communications Research Laboratory (CRL), National  
</dc:description>

<dcterms:extent>xi, 297 p. : ill. ; 30 cm.</dcterms:extent>

<dc:identifier>4924600989</dc:identifier>

<dc:identifier>9784924600980</dc:identifier>

<dc:language>English and Japanese.</dc:language>

FileEditViewInsertFormatToolsDataGo ToFavoritesHelp

Google

Search

Share

Sidewiki

Check

Translate

AutoFill

Page

http://docmgrap02exdu.dev.oclc.org:4321/msd/Crosswalk/Translation\_Tables/dcterms\_marc.xls

http://docmgrap02exdu.dev.oclc.org:4321/msd/Cros...

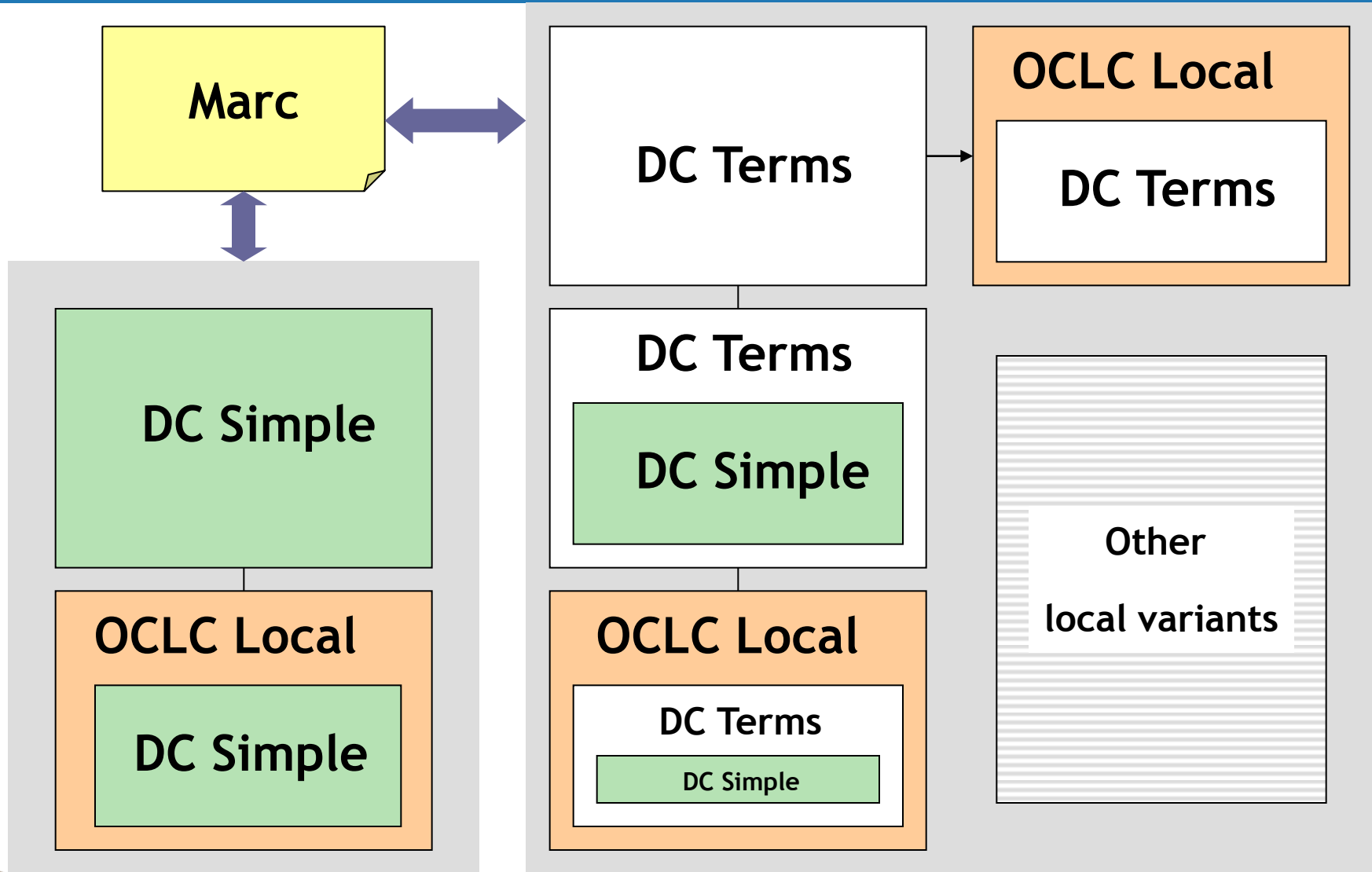
C99fx

	B	C	G	H	I	J	K
16	dc:date		dcterms:Period	DCTERMS.Period	http://purl.org/dc/terms/Period	887	??
1	t-name xml						
2	dcterms:audience						
3	dcterms:educationLe						
4	dcterms:mediator						
5	dc:contributor						
6	dc:coverage						
7	dcterms:spatial						
8	dcterms:spatial						
9	dcterms:spatial						
10	dcterms:spatial						
11	dcterms:spatial						
12	dcterms:temporal						
13	dcterms:temporal						
14	dcterms:temporal						
15	dc:creator						
16	dc:date						
17	dc:date						
18	dc:date						
19	dcterms:available						
20	dcterms:available						
21	dcterms:available						
22	dcterms:created						
23	dcterms:created						
24	dcterms:created						
25	dcterms:dateAccept						
26	dcterms:dateAccept						
27	dcterms:dateAccepted						
28	dcterms:dateCopyrighted		dcterms:Period	DCTERMS.Period	http://purl.org/dc/terms/Period	887	??
29	dcterms:dateCopyrighted		dcterms:W3CDTF	DCTERMS.W3CDTF	http://purl.org/dc/terms/W3CDTF	046	??

DC-MARC crosswalk/

Unknown Zone

# The MARC-Dublin Core relationship



# CDF as an application profile

## Inputs

MARC 21-  
2709

ONIX Books

MARC XML

MODS

DC XML

OAI-DC XML

OCLC CDF

DC-Qualified

ONIX Serials

**OCLC's  
Common  
Data  
Format**

MARC 21-  
2709

## Outputs

OCLC MARC

OCLC CDF

ONIX Books

MARC XML

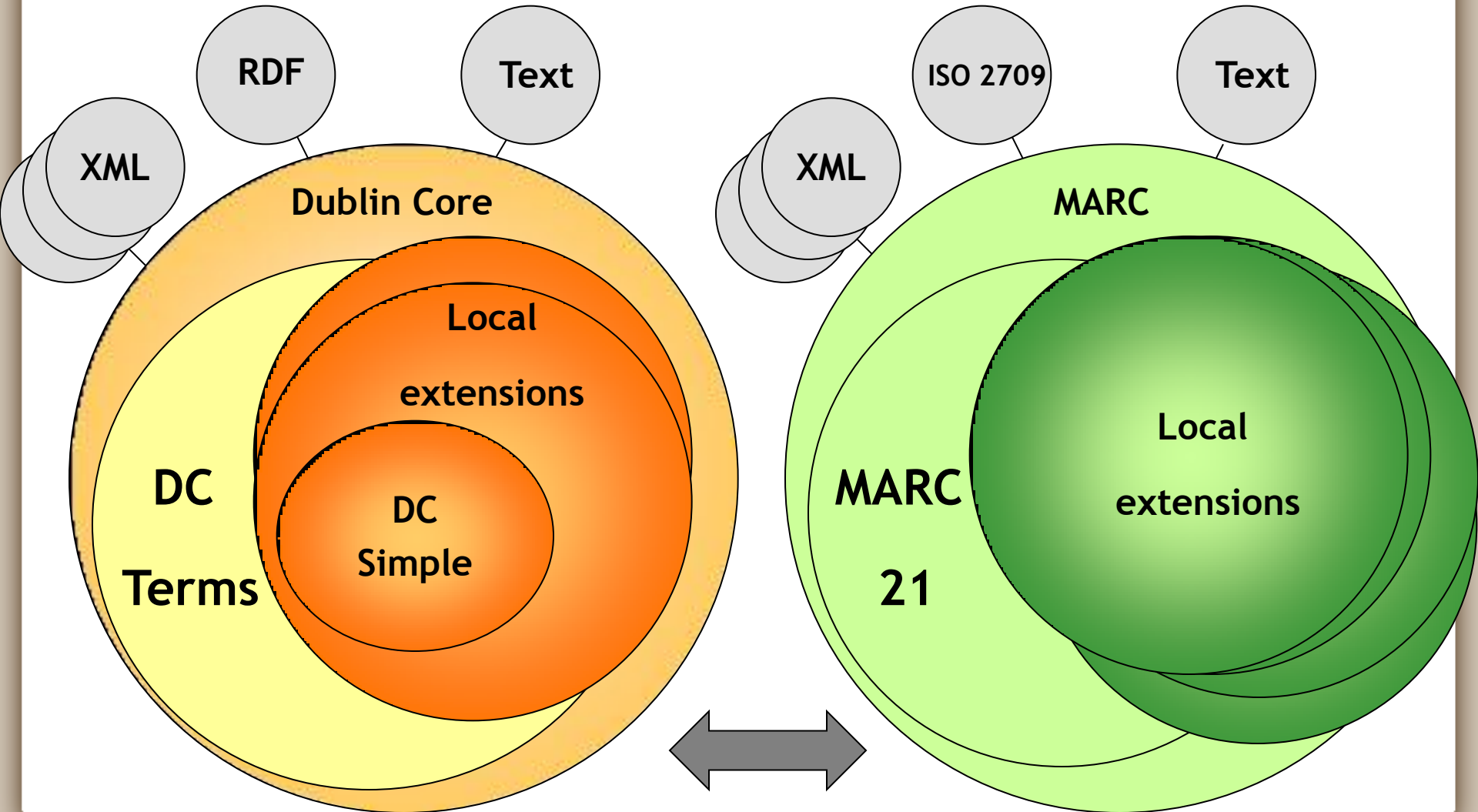
DC XML

DC-Qualified

MODS



# Translations and conversions, expanded





# A graphical user interface to crosswalk maps

## Inputs

<map>  
Source: MARC 245 \$a  
Target: ONIX Title  
</map>

<map>  
Source: MARC 650 \$a  
Target: ONIX Subject  
</map>

Editing  
interface

## Outputs

Search  
interface

Map  
database

Standard  
translation

Implied  
translation

Application  
profile

Version  
upgrade



**In sum:**

**What we have learned to do pretty well**



- Manipulate MARC and MARC-like formats.
- Manage variation in structure and character encoding.
- Manage different versions.
- Reduce and (**almost**) eliminate data processing silos.
- Automate what is systematic and regular.
- Reuse software and metadata subject matter expertise.

# But...

- We must move toward new paradigms for metadata creation and maintenance that permit:
  - Greater interoperability and shared metadata.
  - Mechanisms for allowing metadata to “grow up” over time.



# The old and new paradigms

## Non-MARC elements

Subject

Publisher

Identifier

Contributor

Physical description

## MARC record



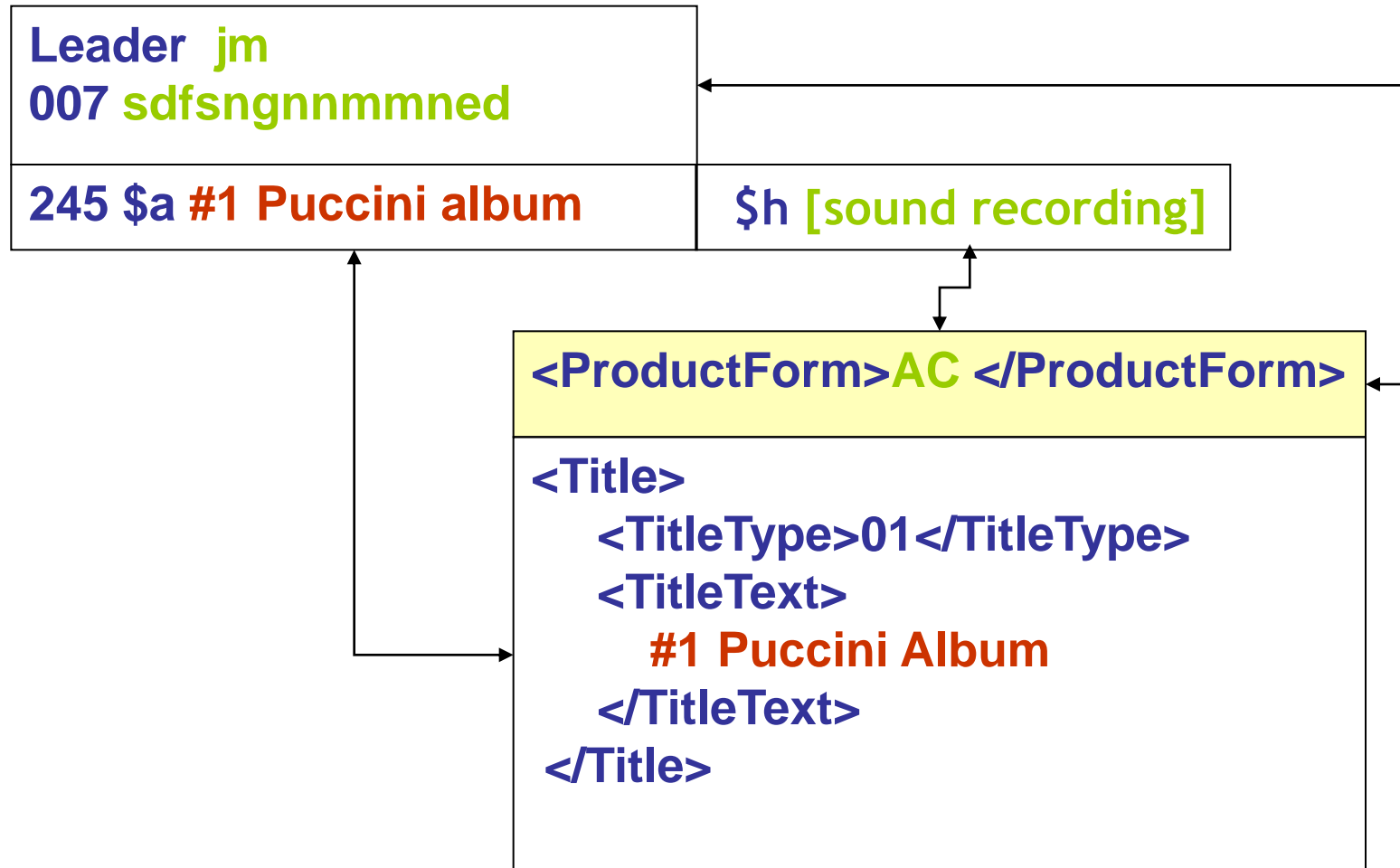
ISBD punctuation

AACR2 encoding

# Some problems

- Extra effort is required to add, validate, and dismantle ISBD and AACR2 rules, adding complexity to the metadata translation/conversion problem space.
- AACR2 introduces concepts that the other standards don't have.
- The ISBD and AACR2 layers are not a worldwide standard.
- Vocabulary and semantic concepts are different.
- Differences in punctuation and formatting require the crosswalk to peek at the data. As a result:
  - The mappings are brittle.
  - Duplicate detection is difficult.

# An example: Physical descriptions in ONIX and MARC



# The two paradigms, redux

## The MARC standard

- Record-oriented
- Tailored to applications in the library community
- Designed for storage
- Static

## Modern non-MARC standards

- Element or field-oriented
- Agnostic about how the data will be used.
- Designed for transmission
- Dynamic

# **Next steps 1: Participate in the transition to RDA**



- **The ONIX/RDA Framework solves some of the problems with physical descriptions by proposing a registered common vocabulary that both standards share.**
- **RDA also has far-reaching recommendations for linking data (instead of copying it).**
- **But the RDA draft standard still has many formatted fields that will be difficult to process algorithmically.**



## **Next steps 2: Apply lessons from studies of MARC field usage**



- Only 21 to 30 tags occur in 10% or more records in OCLC's WorldCat database.
- Fixed-length data elements, identifier fields, main entry and title fields are the most common.
- Only a subset of fields is indexed by library search systems
- Notes fields are common, but machines aren't good at interpreting them.

“**MARC** data cannot continue to exist in its own discrete environment, separate from the rest of the information universe. It will need to be leveraged and used in other domains to reach users in their own networked environments. The 200 or so MARC 21 fields in use must be mapped to simpler schema.”

Source: [Implications of MARC Tag Usage on Library Metadata Practices](#)

# For more information



- Carol Jean Godby. 2010. “Mapping ONIX to MARC.”  
<http://www.oclc.org/research/publications/library/2010/2010-14.pdf>.
- Gordon Dunsire. 2007. “Distinguishing Content from Carrier: The RDA/ONIX framework for Resource Categorization.”  
<http://www.dlib.org/dlib/january07/dunsire/01dunsire.html>.
- EDItEUR. <http://www.editeur.org/>.
- Karen Smith-Yoshimura, Catherine Argus, Timothy J. Dickey, Chew Chiat Naun, Lisa Rowlison de Ortiz, and Hugh Taylor. 2010. “Implications of MARC tag usage on library metadata practices”  
<http://www.oclc.org/research/publications/library/2010/2010-06.pdf>