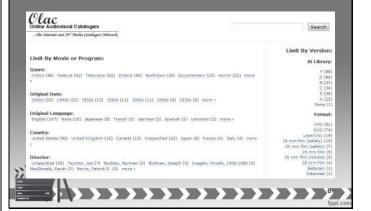
	Mining MADC for Marries Income Date
	Mining MARC for Moving Image Data
Mining MARC for	
<b>.</b>	Mashcat
Moving Image Data	January 13, 2016
Mashcat	Kelley McGrath
January 13, 2016	University of Oregon
	enversity of eregen
Kelley McGrath	
University of Oregon	
	M/bu did Lbooms interacted in the
	Why did I become interested in the
Why?: Input > Output	problem of moving image data in MARC?
	My first job out of library school was as the
Find Videos Clear Form All wideos	a/v cataloger at a university with a large
Search terms are optional.	media collection. I was quickly dismayed by
Original Date: All dates	the disconnect between what I knew when
Action & adventure Epic Science fiction Adaptations Fantasy Shorts	I was cataloging something, what I could
Aviant-garde Historical Sports Biographical Horor Sport	put in the record and what a user could get
Black & white Musical Thriller Children/Family Mystery TV series/movies	back out. This form was my first attempt to
Comedy Romance War Crime Romantic comedy Western	
Setting (Place & Time) & Characters: Region: Anywhere  State (U.S.): Anywhere	bridge that gap. The video collection was in
Country: Anywhere  Time: Any	closed stacks and the videos were only
Types of Characters: Any given accession numbers so there was no	
Captioned Awards: None selected	way to browse the collection. The form
	gives users a few ways to explore the
	collection.
Find Videos Clear Form	All videos 👻
Search for:	as Keyword 👻
Search terms are option	al.
Original Date: All dates - Origin: Any cour	ntry 👻
Genre/Form:	
Action & adventure	Science fiction
Adaptations Fantasy	Shorts
Animation Foreign	Silent
Avant-garde Historical Biographical Horror	Sports Spy
Black & white	Thriller
Children/Family Mystery Comedy Romance	TV series/movies
Comedy Romance Crime Romantic comedy	Western
Setting (Place & Time) & Cha	aracters:
Region: Anywhere	State (U.S.): Anywhere -
Country: Anywhere	Time: Any 👻
Types of Characters: Any	★
Captioned Audio-described Awards: 1	None selected 🗸

Limitations: Complexity Italian comedies from the early 1970s (fiction ADJ (films OR television){655} AND <i>italy</i> {655} AND ((comedy OR comedies)){655}) AND (1970 OR 1971 OR 1972 OR 1973 OR 1974) SAME (motion OR release OR broadcast OR television){500 518})	However, the form has some substantial drawbacks. First, the search strategies are very complex. You can see why a user would never be able to come up with this much less manage to type it in.
<ul> <li>Limitations: Inside Knowledge</li> <li>MARC tags</li> <li>Search options and Boolean operators supported by catalog</li> <li>Local practices</li> </ul>	The search strategies require inside knowledge of the way the catalog software works, of MARC tags and of local cataloging practices.
<ul> <li>Limitations: High Maintenance</li> <li>Retrospective data cleanup</li> <li>Editing of new records to conform to local practices</li> <li>Not transferable or shareable</li> </ul>	The form is also high maintenance. Initially, we did substantial data cleanup to ensure that the data conformed to the expectations of the search strings. In order to make the form work with new records, we had to edit them according to our local practices. Because we relied on local practices, we couldn't effectively share our work.

## Facets: structured data, right data



Facets would be an even better way to enable users to explore the collection. However, you can't exactly facet on "Originally produced as a motion picture in 2006." We need machine-actionable data that is consistently-structured and that answers the questions we care about. The topsy-turvy, convoluted search strategies I used for that early form show that we often don't have that. It isn't practical to start all over again so I started trying to figure out how much we could extract from our existing records.

Online Audiovisual Catalogers	Search
The Internet and ,84' Hedia Catalogers Network	
	Limit By Version:
Limit By Movie or Program:	At Library
Genre: Rotion (96) Feature (62) Television (62) Drama (49) Nonfiction (30) Documentary (23) Horror (22) more	r (56 D (50 8 (37 C (36
Original Date: 2000s (55) 1990s (22) 1950s (13) 1960s (11) 1980s (11) 1890s (8) 1930s (8) more =	E (36 A (22 None (1
Original Language: English (107) None (10) Japanese (8) French (3) German (3) Spenish (2) Unknown (2) more >	Format
Country: United States (90) United Kingdom (15) Canada (13) Unspecified (10) Japan (8) France (4) Italy (4) more	VHS (81 OVD (74 LaserOisc (19 16 mm film (aufaty) (13 35 mm film (aufaty) (7
Director: Unspecified (25) Toynton, Ian (7) Buckley, Norman (3) Bullmen, Joseph (3) Inegeki, Hinzehi, 1903-1980 (3) MacDoneld, Sarah (3) Norma, Patrick R. (3) more +	16 mm film (3 35 mm film (nitrate) (5 35 mm film (4 8etacam (1 Videoreal (1

<ul> <li>Where is the Data?</li> <li>A. In a structured form designed to support search and discovery functions</li> <li>B. In a structured form, but doesn't answer the question</li> <li>C. As free text riddled with typos, often with a wide variety of ways to say the same thing</li> <li>D. All over the place (with no constraints on consistency)</li> <li>E. Nowhere</li> <li>F. All of the above</li> </ul>	When you start looking for the data, it gets complicated very quickly. Sometimes, you have beautiful data, but more often you have data that isn't quite right or that is buried in text strings. It's also common to have too many choices when the same type of information occurs multiple times in a single record. For example, the format of a video can be recorded in several places and there's nothing that enforces consistency. During the transition from VHS to DVD, many catalogers copied a VHS record and edited it to describe the DVD version. Unfortunately, they often forgot to update one of the VHS-related values, which led to a lot of records with an identity crisis—they couldn't make up their minds if they were describing VHS or DVD. Alternatively, sometimes that data just isn't there. The original theatrical release date is important for feature films, but the cataloging rules focus on the publication in hand and don't require this information.
Where is the Data?         Fixed fields: 008/06 +008/07-10 +008/11-14         Psycho (1960)         - p19981960         - s1998         Psycho (1998)         - s1998         Multiple films         End run around MARC: p19601998	Let's look at the various places that the original date of a moving image might appear in a MARC record. Sometimes the original date appears in date2 of the 008, as in the first example. However, this is only done if the video is a straight re- release. If new content of any sort, such as subtitles or special features, has been added, the record is coded with only the date of publication of the video, as in the second example. Sometimes the year of original release and the year the video was published are the same so there's only one date. Since date2 is not repeatable, this method doesn't work for videos that contain more than one film, such as a collection of animated shorts. In another twist, some libraries reverse date1 and date2 so that their public catalog will search and sort by the original date.

Where is the Data? 033 Date of Broadcast: Pertains to the broadcasting (i.e., transmission) or <i>re-broadcasting</i> of sound or visual images. 033 01 \$a19950105	There is a field for date of broadcast, but it isn't limited to the date of the original broadcast.
Where is the Data? Text (headings and notes) 130 True grit (Motion picture : 1969) 500 Originally broadcast on television in 2009. 518 Recorded on Feb. 2, 1991. 505 \$t Tunnel of love / \$r Robert Milton Wallace \$g (1997, b&w, 12 min.)	The original date may appear in different text strings with all the parsing problems that come with that.
Where is the Data? 046 \$k Beginning or single date created 046 \$k 1977 Precise, repeatable "Date or beginning of the date range on which a resource has been created when it is not more appropriately recorded in another field. Dates contained in subfield \$k may not be coded elsewhere in the formats."	Finally, there is a field that would seem to be a good fit except for constraints in the text of the MARC format. OLAC is in the process of trying to get these restrictions removed. No one at the Midwinter MARC Advisory Committee could think of a reason to keep them so I expect our proposal to remove these restrictions will be approved at Annual.

The kitchen sink that is 300\$b Other physical details "Physical characteristics such as illustrative matter, coloration, playing speed, groove characteristics, presence and kind of sound, number of channels, motion picture presentation format, etc."	Going forward, it would be best if there is only one place to record a given type of data. Conversely, only one type of data should be recordable in a given field. Unfortunately, this isn't true of MARC. 300\$b, other physical details, is a particularly egregious example. Nobody in their right mind would define a field like this for data processing.
Bibframe and 300\$b bf:colorContent bf:illustrationNote "sd., col. and b&w ;" http://bibfra.me/vocab/marc/color otherPhysicalDetails "sd., col. and b&w ;"	Both the LC and the Zepheira versions of Bibframe have specific fields for at least some of the data, such as color content, that are currently recorded in 300\$b. Unfortunately, these all appear only to take literal values. This is a step back from MARC, which at least had coded values for color content in 007. In addition, the conversion algorithms that I could find are not very sophisticated and keep everything glommed together in a single textual field.
State       State         St	I looked at a sample of around 90,000 moving image records and almost all instances of 300\$b fall into a narrow range of patterns. I found over 900 variations, although this could be reduced by normalizing the punctuation. Over half of these occur only once. 90% of the 300\$b fields contain one of these strings or are blank so it would be pretty straightforward to write something that would do a better job with conversion to Bibframe.

The long, long tail         s.d., b&w.         b sd., col.,         sd., col. with b&w segments, stereo.         digital, WMV file (1471 Kbps), sd., col.         sd., col. tinted         sd., b&w with tinted and col. sequences         Films for the Humanities & Sciences	There is a long, long tail. There are typos and weird punctuation. There are variant phrasings, such as "segments" instead of the more common "sequences." Although color and sound are the most common types of information in 300\$b for moving images, other information does appear. The cataloging rules do not provide guidance on recording information about tinting and toning, which increases the number of variant forms. And then there are the outliers that are mis-tagged or just make no sense.
<ul> <li>More Specific Data Wanted</li> <li>041 \$h - Language code of original</li> <li>041 \$a - Language code of text/sound track or separate title</li> <li>041 \$j - Language code of subtitles or captions</li> </ul>	With the transition to a new data carrier, we have the opportunity to think about where we might want different or more detailed data. For example, most moving image language data is currently coded in these three subfields.
<ul> <li>More Specific Data Wanted</li> <li>Original language</li> <li>Soundtrack (dubbed or not?)</li> <li>Audio description</li> <li>Intertitles (silent films)</li> <li>Subtitles</li> <li>SDH (Subtitles for the deaf and hard of hearing)</li> <li>Closed-captions</li> <li>Open-captions</li> </ul>	If we really wanted to give users a clear picture of what they're getting, it would be better to have more specific categories.

Consistent, machine-actionable LDR/06 = g, 008/33 = m Run Time: 008/18-20 • 001-999 = # of minutes • = unknown •     = no attempt to code • 000 = over 999 minutes	We currently don't have consistent, machine-actionable data in many places that we need it. Take duration or run time. It can be recorded in the 008, but only three characters are allotted so only times under 1000 minutes can be recorded. Times can only be recorded in one-minute increments, which is not helpful for very short clips.
Consistent, machine-actionable 300\$a • 1 videodisc (120 minutes) • 1 videocassette (1 hr., 34 min., 53 sec.) • 1 videocassette (10, 10, 26 min.) • 2 videocassettes (approximately 60 min. each) • 1 videocassette (ca. 1 hour 30 min. (i.e. 72 min.))	More complex and exact information about duration can be recorded in 300\$a, but with all the drawbacks of text strings.
Consistent, machine-actionable         1 videodisc (120 minutes) <u>Duration</u>	As part of a project I'm working on, I'm trying to convert information from MARC into standardized forms that can be used for analysis and comparison. These five pieces of information enable almost all the variations on duration in my dataset to be normalized in a consistent manner. The time is recorded using a standardized method.

Consistent, machine-actionable1 videocassette (10, 10, 26 min.)DurationDurationDurationDurationDurationDurationP46MtotalcorrectP10Mpart1correctP10Mpart2correctP26Mpart3correct	Oftentimes, the duration of both the whole and the parts is recorded and this can be accounted for.
Consistent, machine-actionable 2 videocassettes	Modifiers such as approximately and over can be added where appropriate.
(approximately 60 min. each)         Duration       Duration       Duration       Duration         Type       Part Number       Qualifier       Validity         P2H       total       approximately       correct         P1H       part       1       approximately       correct         P1H       part       2       approximately       correct	
Duration       Duration <th< td=""><td>In some cases, the cataloger corrects the time stated on the piece and this can also be accounted for. RDA is at least interested in supporting more machine-actionable data and has been investigating how to do this. I am not sure that Bibframe is considering this angle as much as it should.</td></th<>	In some cases, the cataloger corrects the time stated on the piece and this can also be accounted for. RDA is at least interested in supporting more machine-actionable data and has been investigating how to do this. I am not sure that Bibframe is considering this angle as much as it should.

Matching Works         245 \$a Technology         330 \$a Secrets of the superbrands         245 \$a Secrets of the superbrands. \$p Technology.         \$05 \$g disc 1. \$t Technology         \$g disc 2. \$t Fashion \$g disc 3.         \$t Food.	In MARC, the same film or program can be described in more than one way. It would be helpful if all the variations could be identified so that they could be presented to the user as equivalent. OCLC tries to group "works," but, for various reasons, their works are not the same as FRBR works. It is much more difficult to identify all the instances of a true FRBR work. These are three different ways that the Technology episode of Secrets of the Super Brands might appear. The need to manipulate and compare different forms of titles is a strong reason not to reduce titles to flat strings.
Matching Works         245 \$a Friends. \$n 1996-09-26, \$p The one where no one's ready         245 \$a The one where no one's ready         730 \$a Friends (Television program)         245 \$a The one with the Princess Leia fantasy ; \$b The one where no one's ready ;         730 \$a Friends (Television program)         245 \$a The one with the Princess Leia fantasy ; \$b The one where no one's ready ;         730 \$a Friends (Television program)	Here are several ways that the title of a particular episode of Friends might appear.
Matching Works         245 \$a Friends. \$n Season three. \$n Disc one         505 \$a The one where no one's ready         245 \$a Friends. \$n Season three         505 \$a Disc 1 The one where no one's ready         245 \$a Friends. \$p The complete series         505 \$a (Season 3): Disc 9:Ep. 50: The one where no one's ready         245 \$a Friends.         245 \$a The best of Friends         505 \$a v. 2. The one where no one's ready         245 \$a The best of Friends         505 \$a v. 2. The one where no one's ready         730 \$a Friends (Television program)	And here are even more variations.

<section-header></section-header>	Finally, I want to talk about one of my major frustrations with MARC, which I'm not sure that Bibframe is resolving. I think of this as the Humpty Dumpty problem. Many videos contain multiple works, such as a collection of animated shorts. Various things can be said about each short, such as who the director is, when it was made, what language it's in, and so on. These things are usually in different MARC fields. Combining the director of one film with the creation date of another in search results or facets misleads the user.
Humpty Dumpty Problem • \$3 is not for machines • \$8 linking subfield not widely implemented or used in bibs	Sometimes subfield 3 is used to try to connect these pieces of information, but the free text used in \$3 is not very machine friendly. Subfield 8 is designed for this purpose, but I don't know of any systems that implement it for bibliographic records. Nor does \$8 help sort out situations where information about more than one work is combined in a single field.
All the pieces 511 Cyrus Stevens, violin (1st work) ; Pamela Dellal, mezzo-soprano (2nd work) 518 The 1st and 3rd works recorded at the Sonic Temple, Roslindale, MA, Dec. 5 and 14, 2001, respectively 505 Sonata for violin and piano (17:54) A packet for Susan (19:59) 650 Sonatas (Violin and piano) 650 Songs (High voice) with piano 700 Boykan, Martin. Sonatas, violin, piano. 700 Boykan, Martin. Packet for Susan	This is a music example, but the same principle applies to videos. A person looking at a record can usually untangle which pieces of information go with what, but it's hard to imagine how an algorithm could do so.