

FOUNDATIONS of DATA CURATION

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales



School of Information Sciences

• University of Illinois at Urbana-Champaign

A person stands in a dark room, possibly a server room or data center, with rows of server racks filling the background. The person is positioned centrally, facing forward. The lighting is low, with the person's face and the server racks being the primary light sources.

DATA MODELS: TREES

①

TEXT AND DOCUMENTS

Text and Documents

- Why is text important in data curation?
- Examples of data intensive documents.
- The promised functionality of digital documents . . . but not easily realized

What's so important about documents?

The document is the natural unit of textual information.

Why are documents important? Because...

That's where the information is.

Arguably much more information exists in documents and unstructured natural language text than exists in databases.

That's where the action is.

Documents are typically instruments of action; information only has traction on the world when it is communicated in documents

you're hired; you're fired; we agree; you own; you owe;

Even databases typically only have effects when a report (a document) is generated and read by someone (or some processing agent).

That's where we live, work, and play.

We cannot imagine our social lives -- commercial, scientific, cultural everyday -- without the medium of document-based communication.



Documents and data

Relational databases seem a natural fit for certain kinds of data:

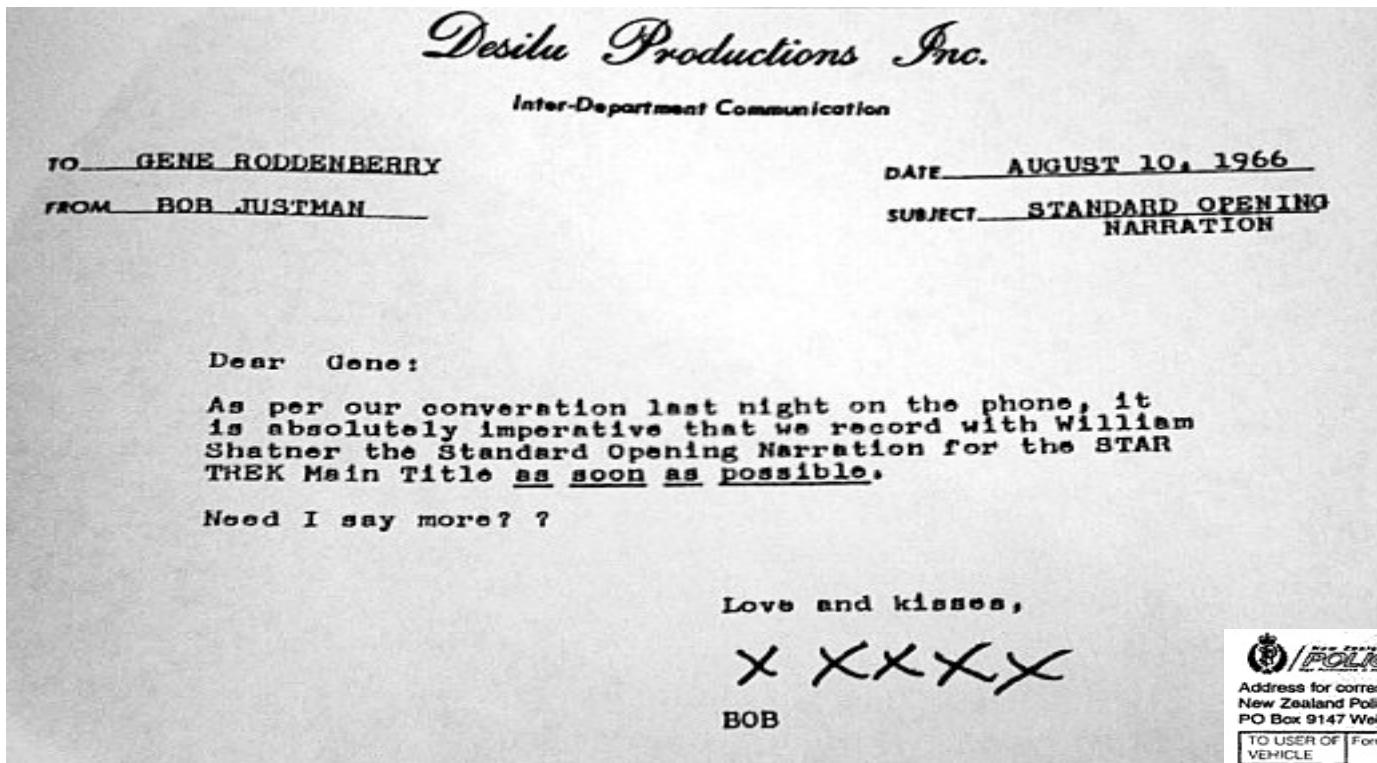
Particularly where the data has this form: *something* has a *value* for an *attribute*

But most the information in the world is contained not in databases, but in the *text of documents*.

This poses challenges for the relational model:

1. In the text of a document is not obvious, and may require considerable human analysis, to see *what* is being said about *what*.
2. The document's text *itself* often needs to be organized and managed (as in publishing applications), rather than the data being asserted by the text; most documents do not appear to be tabular in nature.

Examples



 POLICE

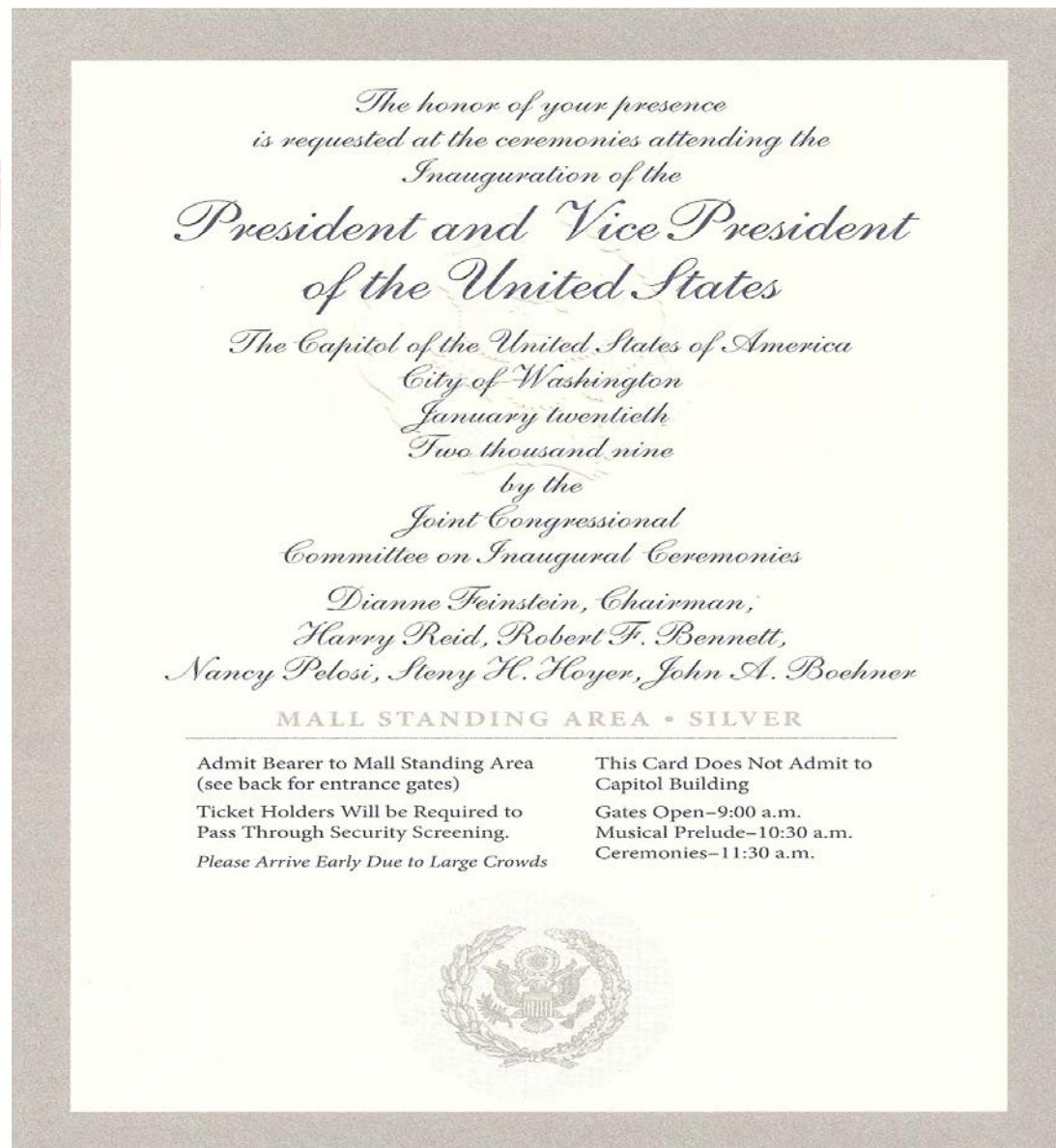
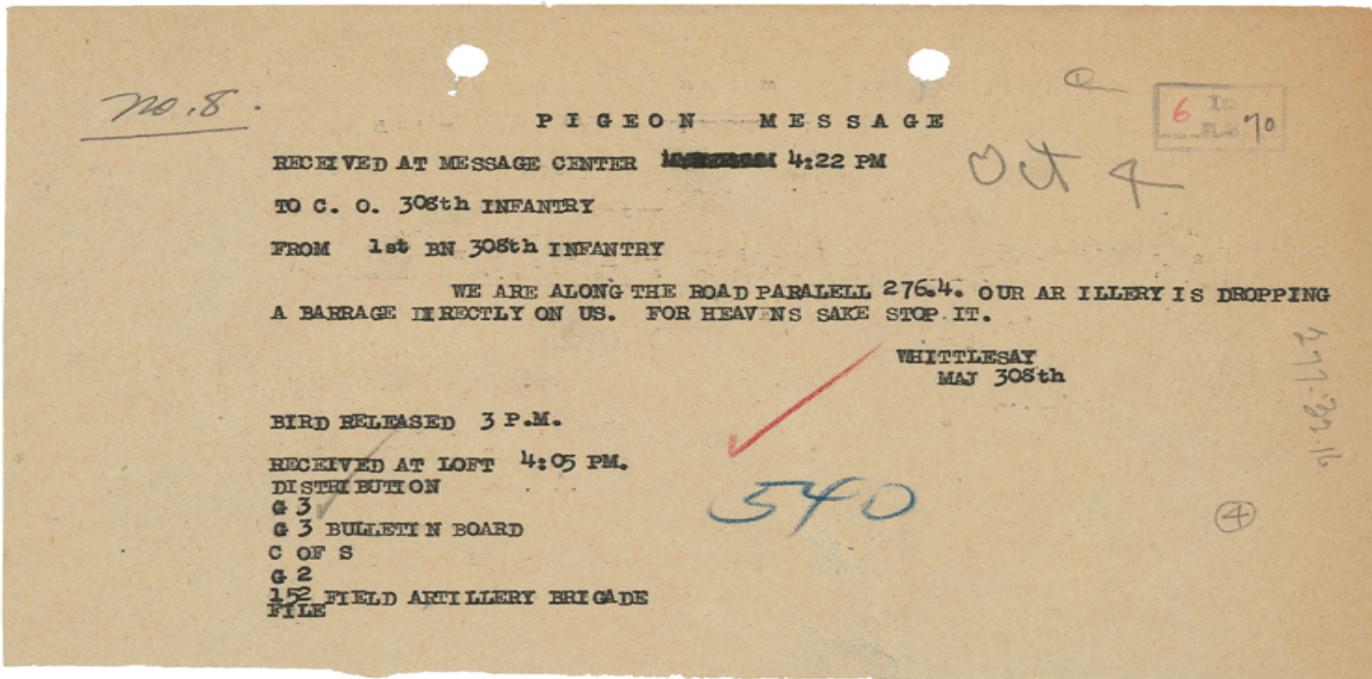
INFRINGEMENT NOTICE
(ISSUED UNDER THE AUTHORITY OF THE LAND TRANSPORT ACT 1968)

POL 405
9/2002

NOTICE NUMBER N 3735700

TO USER OF VEHICLE	Forename(s) <i>Justin Alexander</i>	Family Name <i>LEE</i>
Address <i>[Redacted]</i>		
Occupation <i>Accountant</i>	Date of Birth <i>23/6/1974</i>	Driver Licence Number <i>BP 086926</i>
ALLEGED INFRINGEMENT OFFENCE(S) DETAILS		
Date of Offence <i>23/6/1974</i>	Time <i>18:25</i> (24 Hour Clock)	Day of Week <i>S M T W T F S</i>
Vehicle Type <i>Sedan</i>	Vehicle Make <i>Honda</i>	Reg. No. <i>AEH 924</i>
Road/Street <i>STATE Hwy/Way one</i>	Locality <i>POKERU</i>	Infringement Fee Payable
Offence Number <i>1</i>	Offence <i>Exceeded 100 mph</i>	\$ 120 -

Examples



Examples

Tele: "DAILY SERVICE" Phone: 32.

Nepal Transport Service		
BUS NO.	Class	DATE
FROM AMLEKHGUNJ		TO KATHMANDU
S.	Rs.	Sd.
PASSENGERS TICKET		

CERTIFIED COPY OF AN ENTRY OF MARRIAGE

GIVEN AT THE GENERAL REGISTER OFFICE

Application Number: COL Number

19 Year		Marriage solemnized at		The Register Office		in the	
District of		County Name		in the		Borough Name	
No.	When married	Name and surname	Age	Condition	Rank or profession	Residence at the time of marriage	Father's name and surname
No.	Date	Name and Surname	Age	Condition	Job Title	Address	Father's Name and Surname
	Month						Rank or profession of father
	Year	Name and Surname	Age	Condition	Job Title	Address	Father's Profession
Married in the		Register Office				by	Register Name
This marriage was solemnized between us,		Groom Signature	In the presence of us,	Witness Signature 1		Register Address	
		Bride Signature		Witness Signature 2			

SAMPLE CERTIFICATE

CERTIFIED to be a true copy of an entry in the certified copy of a register of Marriages in the Registration District of _____ Given at the GENERAL REGISTER OFFICE, under the Seal of the said Office, the _____ day of _____

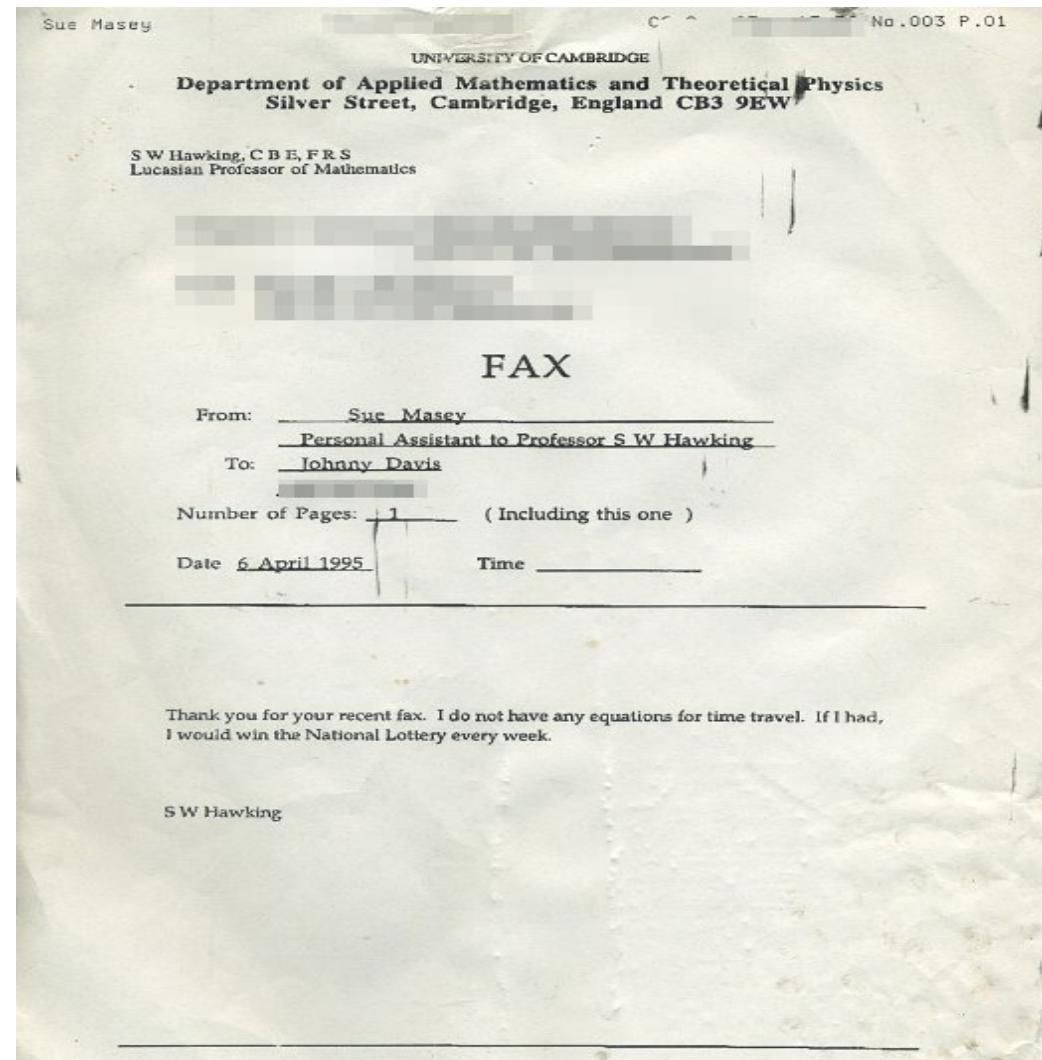
This certificate is issued in pursuance of section 65 of the Marriage Act 1949. Sub-section 5 of that section provides that any certified copy of an entry purporting to be sealed or stamped with the seal of the General Register Office shall be received as evidence of the marriage to which it relates without any further or other proof of the entry, and no certified copy purporting to have been given in the said Office shall be of any force or effect unless it is sealed or stamped as aforesaid.

CAUTION: THERE ARE OFFENCES RELATING TO FALSIFYING OR ALTERING A CERTIFICATE AND USING OR POSSESSING A FALSE CERTIFICATE. ©CROWN COPYRIGHT

MXD000000

WARNING: A CERTIFICATE IS NOT EVIDENCE OF IDENTITY.





is done during the "Reynolds averaging" process, corresponds to covariance terms such as $\langle C_1 \cdot C_2 \rangle$ for a reaction involving two gas phase components 1 and 2. The second-order equation describing the flux term takes the form:

$$\frac{d(C_1 \cdot C_2)}{dz} = -w^2 \frac{\partial^2 C_1}{\partial z^2} - \frac{\partial w \cdot C_1 \cdot C_2}{\partial z} \quad (2)$$

$$+ \frac{2}{3} \cdot \frac{\partial C_1}{\partial z} \cdot \frac{\partial C_2}{\partial z} + \frac{1}{3} \frac{\partial^2 C_1}{\partial z^2} \cdot C_2 + \epsilon \cdot w \cdot \partial_z C_1 \cdot C_2 \quad (3)$$

Flux term (1) is the diffusional transport term, (2) is the vertical gradient of the turbulent flux, (3) is the change due to buoyancy (Eq. (4) is known as the pressure dissipation term and represents corrections of pressure derivatives with concentrations C_1 and C_2). The term ϵ represents the effect of change in flux due to chemical reactions.

The dynamic evolution of the 1-D PBL is described in the current model by using an algebraic stress model (ASM) assumption to close the second-order transport terms. This approach yields equations that would conform to a level-2 closure of the Reynolds averaged equations. We adopted a similar approach (ASME-like) and assumptions to close the second-order chemical equation. This yields the following equation for the flux term:

$$\frac{d(C_1 \cdot C_2)}{dz} = \frac{1}{2} \left(\frac{\partial C_1}{\partial z} \right)^2 + \frac{\partial^2 C_1}{\partial z^2} - \epsilon \cdot w \cdot \partial_z C_1 \cdot C_2 - \frac{P_{AG}}{w} \quad (4)$$

ϵ is the respective correlation term and is defined as (Eq. (5)). As is the usual practice in ASM-based modeling, the flux terms can now be represented in a form similar to the eddy mixing coefficient formulation, for example as follows:

$$\epsilon \cdot w \cdot \partial_z C_1 \cdot C_2 = \frac{2k}{3} \cdot w^2 \cdot \frac{\partial^2 C_1}{\partial z^2} \quad (5)$$

Here k is a constant derived from reanalysis (2). This constant is dependent on the flow conditions and chemical species considered.

The Transport Terms in the chemical transport equations are solved by using a fully implicit finite difference scheme. The vertical grid employed has 40 levels and put on a logarithmic axis to give highest resolution in the lowest 50 m of the model. The model is solved with time steps of 5 s for the chemistry and 6 s for the dynamics equations. The primary focus of these calculations is on evaluating the effect of chemistry on the calculated flux of NO from soil.

PRELIMINARY RESULTS

In the initial set of calculations, the moving coefficients derived as described by (6) were employed in applying the transport equation for NO, NO_2 , and O_3 . Applying the eddy mixing coefficients for the rest of the trace gases, same as that calculated for the temperature in the model. This set of calculations is referred to as "with reaction." In a second set of calculations, the eddy mixing coefficients for all of the trace gases in the model are set to those derived for the temperature and are referred to as "no reaction" in the following discussions. In both of these calculations, the

Technical Support Documentation Page		
1. Expert No.: FHWA-HET-12-046	2. Government Accession No.:	3. Recipient's Catalog No.:
4. Title and Date:	5. Report Date:	6. Recipient Organization Code:
Asset Sustainability Index: A Proposed Measure for Long-Term Performance	July 2012	7. Recipient Organization Report No.:
7. Author(s): Gordon Proctor, Shreya Varma, Steve Varnadoz	8. Recipient Organization Name:	9. Recipient Organization Address:
Gordon Proctor & Associates, Inc.	Starline Corp.	3157 Woodstock Drive Lewis Center, Ohio 43035
Dublin, Ohio 43016	National Center for Pavement Preservation	7457 July Road Okemos, MI 48864
10. Sponsoring Agency Name and Address:	11. Type of report and period covered:	12. Sponsoring Agency Name and Address:
FHWA Surface Transportation Environment and Planning Cooperative Research Program, and the Office of Asset Management, Pavements and Construction	2012	FHWA Surface Transportation Environment and Planning Cooperative Research Program, and the Office of Asset Management, Pavements and Construction
13. Supplementary Notes:	14. Sponsoring Agency Code:	15. Key Words:
This report examines the concept of asset sustainability metrics. Such metrics address the long-term performance of highway assets based upon expected expenditure levels. It examines how such metrics are used in Australia, Britain and the private sector. It also reviews asset management data from selected states to illustrate that long-term sustainability metrics could be produced using available US asset management data.	No restrictions. This document is available to the public from the FHWA Surface Transportation Environment and Planning Cooperative Research Program and the Office of Asset Management, Pavements and Construction. naca.dot.gov/nca/infrastructure/assessment	Asset Sustainability, Asset Management, Long-term Performance, Sustainable Infrastructure, Performance Management
16. Security Classification of this report:	17. Security Classification of this report:	18. No longer:
Unclassified	Unclassified	110
19. Price:	20. Price:	Free

Form DOT F-100-10-72 Reproduction of copyrighted page authorized.



NEUTAJOVANÉ

MS-0144G-50420-00

Slaňovací rám SR-1

Provozní kontroly

	Strana
Slaňovací rám SR-1 - Provozní kontroly	1

	Strana
tabulek	1

	Strana
1. Odkazy	1

	Strana
2. Požadavky na pracovní silu	1

	Strana
3. Požadavek na pracovní síly	1

	Strana
4. Kontrola ovládání sláňovacího rámu	2

	Strana
5. Kontrola ovládání haku TYLER 442 s otevírací pákou T-Handle Release Unit	3

Odkazy

Tab. 1 Effectively axiomatized theories

- Given a derivation of the sentence φ from the axioms of the theory T using the background logical proof system, we will say that φ is a theorem of T . Using the standard abbreviatory symbol, we write: $T \vdash \varphi$.
- A theory T is sound iff every theorem of T is true (i.e., true on the interpretation built into T 's language). Soundness is, of course, normally a matter of having true axioms and a truth-preserving proof system.
- A theory T is effectively decidable iff the property of being a theorem of T is an effectively decidable property – i.e., iff there is an algorithmic procedure for determining, for any given sentence φ of T 's language, whether or not $T \vdash \varphi$.
- Assume now that T has a standard negation connective \neg . A theory T decides the sentence φ iff either $T \vdash \varphi$ or $T \vdash \neg\varphi$. A theory T correctly decides φ just when, if φ is true (on the interpretation built into T 's language), $T \vdash \varphi$, and if φ is false, $T \vdash \neg\varphi$.
- A theory T is negation-complete iff T decides every sentence φ of its language (i.e., for every sentence φ , either $T \vdash \varphi$ or $T \vdash \neg\varphi$).
- T is inconsistent iff for some sentence φ , we have both $T \vdash \varphi$ and $T \vdash \neg\varphi$.

Note our decision to restrict the theorems, properly so called, to the derivable sentences; we will, with free variables derived as we go along through a proof, don't count. This decision is for convenience as much as anything, and nothing hangs on it.

Here's a very elementary toy example to illustrate some of these definitions. Consider a trivial pair of theories, T_1 and T_2 , whose shared language consists of the (interpreted) propositional atoms ' p ', ' q ', ' r ' together with all the wffs that can be constructed from them, using the familiar propositional connectives, whose shared underlying logic is a standard natural deduction system for propositional logic, and whose axioms are respectively:

$$T_1 \vdash p,$$

$$T_2 \vdash (q \wedge r) \rightarrow p.$$

T_1 and T_2 are then both axiomatized formal theories. For it is effectively decidable what is a wff of the theory, and whether a purported proof is a proof from the given axioms. Both theories are consistent. Moreover, both are decidable theories: just use the truth-table test to determine whether a candidate theorem really follows from the axioms.

However, note that although T_1 is a decidable theory that doesn't mean T_1 decides every wff: it doesn't decide e.g. the wff ' $(q \wedge r)$ ', since T_1 's sole axiom doesn't entail either ' $(q \wedge r)$ ' or ' $\neg(q \wedge r)$ '. To stress the point: it is one thing to have an algorithm for deciding what is a theorem; it is another thing for a

MS-0144G-50420-00

NEUTAJOVANÉ



Postup

Provozní kontrola funkčnosti sláňovacího rámu SR-1

- Upřídit se, že SR (Obr. 1 [1]) je v transportní poloze.

- Vytáhnout držadlo (Obr. 1 [2]) a rukou umístit SR (Obr. 1 [1]) do sklopné polohy.

- Uvázit držadlo (Obr. 1 [2]) se musí automaticky zavřít ve sklopné poloze SR.

- Upřídit se, že SR (Obr. 1 [1]) je zapojen ve sklopné poloze.

- Zvednout konzolu závěsu (Obr. 1 [3]), demontujte čep kluzáku (Obr. 1 [4]) a umístit konzolu závěsu (Obr. 1 [3]) do specifické polohy a čep kluzáku (Obr. 1 [4]) do držáku (Obr. 1 [5]).

- SR (Obr. 1 [1]) je v pracovní poloze.

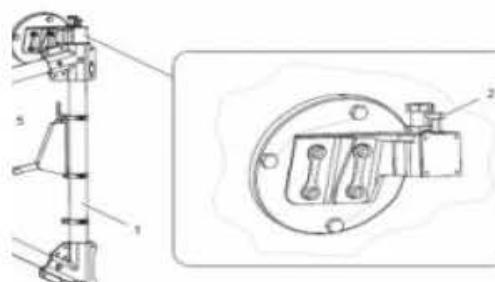
- Podložit konzolu závěsu (Obr. 1 [3]), demontujte čep kluzáku (Obr. 1 [4]) a umístit konzolu závěsu (Obr. 1 [3]) do specifické polohy a čep kluzáku (Obr. 1 [4]) do držáku (Obr. 1 [5]).

- SR (Obr. 1 [1]) je v sklopné poloze.

- Uvázit držadlo (Obr. 1 [2]) a rukou umístit SR (Obr. 1 [1]) do transportní polohy.

- Uvázit držadlo (Obr. 1 [2]) se musí automaticky zavřít v transportní poloze SR.

- Upřídit se, že SR (Obr. 1 [1]) je zapojen v transportní poloze.



Obr. 1 Rám

(A-A)

Rám

ROZS. 000-01

(NTI) using identical construction source as distinct as mine other. Since offers a proof that NTI follows from L. (D-693, p. 311). Agreed with these principles, Shreya argues as follows. C-95C and C-695C had the same characteristics as described by PE, they were identical in that year. So we have:

$$(D-693, C-95C = D-695C)$$

From (D) and (C) we want deduce that C-95C and C-695C are identical principles.

$$(C-95C = D-695C)$$

Now the conclusion about C-695C was:

(D) C-695C never changes in membership, and this in combination with (C) yields

$$(C-95C never changes in membership)$$

That is, as oppositely variable class cannot vary after all.

$$Q. \bar{Q} \in D$$

This is fine as far as it goes, but we're left to wonder why precisely (D) is true. So I think it's true, but what I don't know is the question we want answer? Shreya even says that on all classes have to be invariant if we take L, but why must classes such as C-95C be invariant in the first place?

A similar problem would arise for trying to make Shreya's argument work. The problem is that the notion of "axiom" is not well-defined in her context. The notion of "axiom" is usually understood as something that cannot be derived in another theory. That is, an axiom is a proposition that cannot be derived from any other proposition in another theory. This is the case for the axioms of logic, for example. In the case of NTI, the axioms of my number theory (NTI), things logical in my positive world cannot be derived in another. And in place of (D) we would use the assumption that there are some axioms that can never be derived in any positive world. But what is the ground for this assumption? This is the very thing we want explained.

Whereas the answer, it will have to invoke principles in addition (PE and NG). To see this, consider another \bar{Q}_1 , the identity of individuals:

$$(D) \text{Nexity: } (\bar{Q}_1 = \bar{Q}_2 \rightarrow \bar{Q}_1 = \bar{Q}_2)$$

This is equivalent to PE, so if PE and NG implied that a set has all its members exactly, ID and NG single Element is imply that a class has all its properties exactly. But of course there is no such implication.

The promise of digital documents

This is the grand old dream of radical new functionality.

(cf Paul Otlet, Vannevar Bush, Douglas Engelbart, and Ted Nelson)

- computationally available data items accessible with discipline-specific tools (chemical formulae, proteins, equations, etc.)
 - advanced navigation and viewing optimized for domain-specific browsing and analysis,
 - typed hypertext linking with links as first class objects,
 - data-driven interactive diagrams and graphics
 - computable equations,
 - supportive ontological inferencing
 - thoroughgoing interoperability with other tools
- ... and so on, and on, and on.

Are we achieving the promise of digital documents?

We are not.

What the problems are,

and why we have these problems

is the topic of the next video.

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.

Comments and corrections to: renear@illinois.edu.

The Problem (it's the same problem)

- The situation (circa 1960)

Text is stored and processed in radically different ways

Interaction with text is immediately and directly via storage and processing methods

Explicit and formal conceptualization of text components as such is rare
(and typically only in human memory)

- Why is this a problem?

Huge operational inefficiencies

Lack of functionality

Lack of data independence

The promise of digital documents – unfulfilled

In the last video we noted the importance of documents and the long-standing promise of digital documents to provide exciting new functionality.

But we also noted that we are only part way there

Creating complex, high-performance documents remains arduous and the results are hardly the sophisticated high-performance information environments we were promised.

The Problem (*it's the same problem!!*)

There are many many ways to represent text in documents.

.ll 3i

.mk a

.ce

Preamble

.sp

We, the people of the United States, in order to form a more perfect Union

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	ANSI ASCII
00000000	D0	CF	11	E0	A1	B1	1A	E1	00	00	00	00	00	00	00	00	ĐÍ à;± á
00000010	00	00	00	00	00	00	00	00	3B	00	03	00	FE	FF	09	00	;
00000020	06	00	00	00	00	00	00	00	00	00	00	01	00	00	00	þý	
00000030	11	00	00	00	00	00	00	00	00	10	00	00	02	00	00	00	
00000040	01	00	00	00	FE	FF	FF	FF	00	00	00	00	00	00	00	þýýý	
00000050	FF	ÿÿÿÿÿÿÿÿÿÿÿÿÿÿ															
00001A00	57	00	65	00	20	00	74	00	68	00	65	00	20	00	50	00	We the P
00001A10	65	00	6F	00	70	00	6C	00	65	00	20	00	6F	00	66	00	e o p l e o f
00001A20	20	00	74	00	68	00	65	00	20	00	55	00	6E	00	69	00	t h e U n i
00001A30	74	00	65	00	64	00	20	00	53	00	74	00	61	00	74	00	t e d S t a t
00001A40	65	00	73	00	2C	00	20	00	69	00	6E	00	20	00	4F	00	e s , i n O
00001A50	72	00	64	00	65	00	72	00	20	00	74	00	6F	00	20	00	r d e r t o
00001A60	66	00	6F	00	72	00	6D	00	20	00	61	00	20	00	6D	00	f o r m a m

```
{\rtf1\ansi{\fonttbl{\f0\fswiss Helvetica;}}\f0\pard  
This is some {\b bold} text.\par }
```

```
<w:body>  
<w:p w:rsidR="00146B24" w:rsidRDefault="00146B24">  
<w:bookmarkStart w:id="0" w:name="_GoBack"/>  
<w:bookmarkEnd w:id="0"/>  
</w:p>  
<w:p w:rsidR="00146B24" w:rsidRDefault="00146B24" w:rsidP="00146B24">  
<w:pPr>  
<w:pStyle w:val="Heading1"/>  
</w:pPr>  
<w:r>  
<w:t>Preamble</w:t>  
</w:r>  
</w:p>  
<w:p w:rsidR="001C180C" w:rsidRDefault="00146B24">  
<w:r>  
<w:t xml:space="preserve">We the People of the United States, </w:t>  
</w:r>  
...
```

Interaction is typically directly with these storage structures

The fundamental principles of abstraction and indirection are not implemented



The problems that result

- Training does not transfer
- Tools are not interoperable
- Data from multiple sources cannot be integrated
- Applications development is arduous
- Documentation is difficult
- Validation and assurance is difficult
- Specialized applications (searching, analysis, etc.) are not supported
- Schemas are typically non-existent, or unhelpful

And so on...

Sound familiar?

A closer look at one example: electronic publishing in the 1960s

Input file contains this data:

```
.pa odd;.font Times;.size 14;  
.it;.ce;.in +5 -5;.sk 2p a;.kp next;.toc include; The Sick  
Rose[...]
```

After processing the output is rendered like this:

The Sick Rose
[...]

What are some problems with this approach to organizing text?

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.

Comments and corrections to: renear@illinois.edu.

③

THE SOLUTION: 1. DESCRIPTIVE MARKUP

The Solution: (1) Descriptive Markup (it's the same solution)

The problem, again

The solution: (1) Descriptive Markup

- How it emerged

- How it works and delivers the goods

- Why it works

Our example:

Recall athe problems we noticed with this approach to organizing text.

Input file contains this data:

```
.pa odd;.font Times;.size 14;  
.it;.ce;.in +5 -5;.sk 2p a;.kp next;.toc include; The Sick  
Rose[...]
```

After processing the output is rendered like this:

The Sick Rose

[...]

The first improvement

(from the US Government Printing Office in the 1960s)

A macro is defined to abbreviate formatting commands:

```
&format17 =df
  ".pa odd;.font Times;.size 14;.it;.ce;.in +5 -5;.sk 3p;.sk 2p
  a;.kp next;.toc include;"
```

The macro is used in the input:

This helps a bit (why?)

```
.format17:The Sick Rose [...]
```

But not as much as it might help (why?).

A Problem

Although `format17` helps, it doesn't go far enough.

Since `format17` abstracts to a typographic "look" it can be used wherever you want that look: titles, captions, extract labels, ...

<code>.format17:The Sick Rose</code>	<code>.format17:Fig. 1. A tea rose</code>
<p><i>The Sick Rose</i></p> <p>[...]</p>	 <p><i>Fig. 1. A tea rose</i></p>

So what might be even better?

Don't identify the *look*, identify the component itself

A much better improvement

A macro is defined to identify the logical component of the text itself not the the intended processing, or the appearance of the that component

```
&title. =df
        ".pa odd;.font Times;.size 14;.it;.ce;.in +5 -5;.sk
3p;.sk 2p a;.kp next;.toc include;"
```

The macro is used in the input file like this:

```
.title:The Sick Rose [...]
```

Abstraction from storage

Consider this text:

An example of the Tea rose is: <http://www.example.org/rose/hybrid/tea/pink/42>

It's a mouthful. And what if the image moves?

A entity name is defined to abbreviate a location:

&rose42 =df <http://www.example.org/rose/hybrid/tea/pink/42>

The entity name is used in the input:

An example of the Tea rose is: &rose42;

After processing the output is rendered like this:

An example of the Tea rose is:



Again: Abstraction and Indirection

Again our solution is abstraction

In explicitly identifying recurring logical objects (like titles) we abstract away from the varied and varying details of processing and storage

We then exploit this abstraction by using indirection:

- Mapping object instances to storage locations
- Mapping types of objects to processing rules

Achieving efficiencies and new functionality

Examples of text components

- Title
- Author
- Date
- Abstract
- Section, subsection, subsubsection
- Section title, subsection title ... etc
- Paragraph
- Extract (long quotation)
- Equation
- Diagram
- Footnote

Genre-specific text components

Scientific article:

- Title, author, affiliation, address, date submitted, date revised, keywords, abstract, introduction, methodology, results, discussion, conclusion, diagram, equation, plate, graph, chart, bibliography, bibliography item, date

Playscripts:

- Act, scene, stage direction, line, character, cast list

Poetry:

- Title, author, verse, stanza, couplet, line, half-line

Also:

- Legal and financial documents such as contracts, deeds, licenses, writs, tickets, receipts
- Office documents such as project proposals, monthly reports, position descriptions, performance evaluations, other forms
- Etc...

An XML example

```
<anthology>
  <poem>
    <heading>THE SICK ROSE</heading>
    <stanza>
      <line>O Rose thou art sick.</line>
      <line>The invisible worm,</line>
      <line>That flies in the night</line>
      <line>In the howling storm:</line>
    </stanza>
    <stanza>
      <line>Has found out thy bed</line>
      <line>Of crimson joy:</line>
      <line>And his dark secret love</line>
      <line>Does thy life destroy.</line>
    </stanza>
  </poem>
  <!-- more poems go here -->
</anthology>
```

**Example from the Text Encoding Initiative (P5)*

Other examples

NeXML for phylogenetic information

```
<characters otus="tax1" id="m1" xsi:type="nex:DnaSeqs">
  <!-- ... -->
  <matrix aligned="1">
    <row id="r1" otu="t1"><seq>AACATATCTC</seq> </row>
    <row id="r2" otu="t2"><seq>ATACCAGCAT</seq> </row>
    <row id="r3" otu="t3"><seq>GAGGGTATGG</seq> </row>
    <row id="r4" otu="t4"><seq>GGTCTTAGAG</seq> </row>
    <row id="r5" otu="t5"><seq>CGTCACAGTG</seq> </row>
  </matrix>
</characters>
```

Medical Markup Language

```
...
<clinical_document_header>
<document_type_cd DN="MML Document" S="1.2.392.114319.1.1" V="0300" />
  <provider>...</provider>
  <patient>...</patient>
<!-- ...-->
</clinical_document_header>
<body>
  <!-- ...-->
  <mml:docInfo contentModuleType="report"
    moduleVersion="http://www.medxml.net/MML/report/1.0">
    <mml:securityLevel>...</mml:securityLevel>
    <mml:title generationPurpose="reportRadiology">CT scan Report</mml:title>
    <mml:docID><mml:uid>JPN432101234567RR20020823_CT_20020851501</mml:uid></mml:docID>
  </mml:docInfo>
</body>
...

```

Music Markup Language

```
<section>
  <measure n="0" xml:id="m0" type="upbeat">
    <staff n="1">
      <layer n="1">
        <rest dur="4"/>
      </layer>
    </staff>
    <staff n="2">
      <layer n="1">
        <beam>
          <note xml:id="m0_s2_e1" pname="e" oct="5" dur="8" stem.dir="down"/>
          <note xml:id="m0_s2_e2" pname="f" oct="5" dur="8" stem.dir="down"/>
        </beam>
      </layer>
    </staff>
  ...
</section>
```

Descriptive Markup

Descriptive markup describes the *logical components* of documents.

It does not specify *processing*.

Advantages of descriptive markup

Authoring, Editing, Transcribing:

- Composition is simplified
- Writing tools are supported
- Alternative views and links facilitated

Publishing:

- Formatting generically specified and modified
- Apparatus automated
- Output device support enhanced
- Portability maximized

Retrieval and Analysis:

- Information retrieval supported
- Analytical procedures supported

How does this help with data curation

Descriptive markup makes digital documents...

- easier to create
- easier to maintain
- easier to convert (new formats, new delivery software)
- better integrated with workflow in organization
- better integrated with other applications and tools (databases, word processing templates, indexes,)etc.
- more accessible to varied audiences
- easier to accommodate different technological circumstances (varying hardware, operating systems, browser software (brands and versions both), connectivity (bandwidth), etc.
- Easier to accommodate different perceptual abilities (blindness, other sight disabilities, dyslexia, etc.)

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

**Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.**

Comments and corrections to: renear@illinois.edu.



FOUNDATIONS of DATA CURATION

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales



School of Information Sciences

• University of Illinois at Urbana-Champaign



DATA MODELS: TREES

④

THE SOLUTION: 2. TREES

The Solution: (2) Trees

- The OHCO model of text
- The data structure here is a tree, a kind of graph
- Trees can be serialized in formal languages defined by context free grammars

The OHCO model of text emerges

Text is an Ordered Hierarchy of Content Objects

- *content objects* = things such as chapters, paragraphs, sentences, stanzas, lines, speeches, equations, titles, headings, abstracts
- *hierarchy* = sentences inside paragraphs, paragraphs inside sections, sections inside chapters, etc., nesting with no overlaps
- *ordered* = objects proceed or follow one another

Trees

Most documents can be modeled as *trees*.

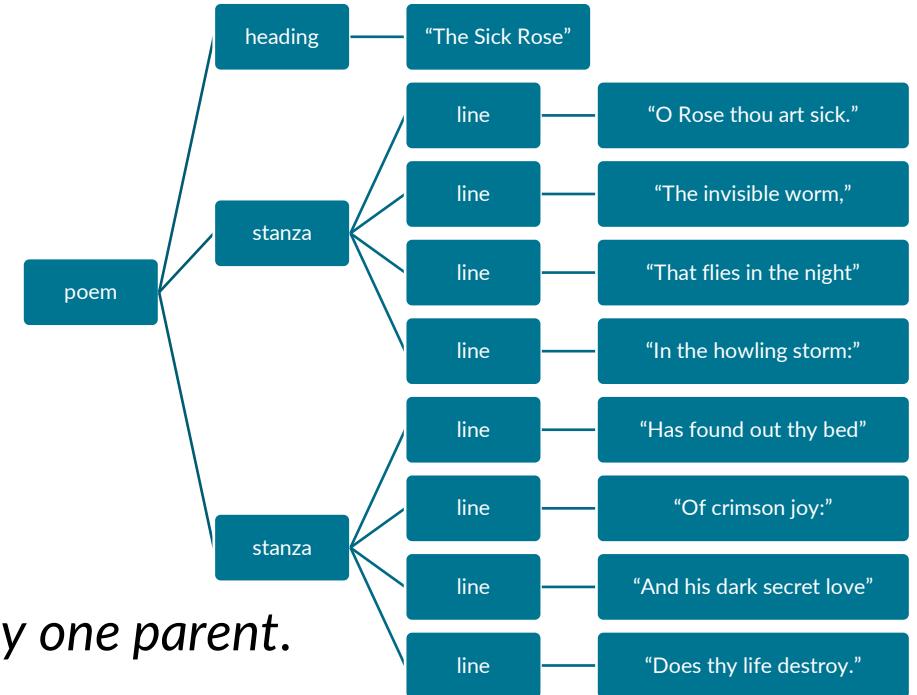
A tree, in our sense*, is

*a directed acyclic graph with ordered branches
and all nodes except one having exactly one parent.*

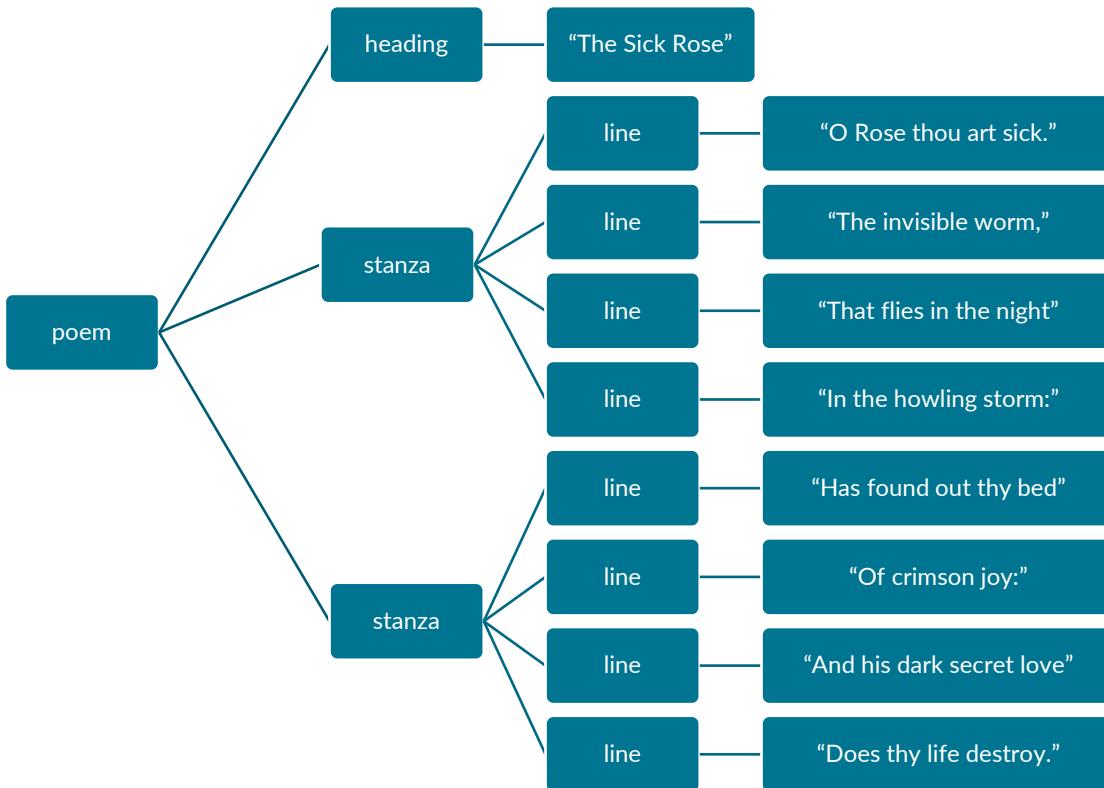
The nodes are labelled (e.g. *title*)

and, also, typically annotated with additional information (such as lang=English).

* A specialization of the usual definition.



Using XML to serialize a tree



```
<poem>
  <heading>The SICK ROSE</heading>
  <stanza>
    <line>O Rose thou art sick.</line>
    <line>The invisible worm,</line>
    <line>That flies in the night</line>
    <line>In the howling storm:</line>
  </stanza>
  <stanza>
    <line>Has found out thy bed</line>
    <line>Of crimson joy:</line>
    <line>And his dark secret love</line>
    <line>Does thy life destroy.</line>
  </stanza>
</poem>
```

A tree can be serialized with a formal language defined by a context free grammar, such as an XML language.

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

**Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.**

Comments and corrections to: renear@illinois.edu.

⑤

WHY THE SOLUTION WORKS

Clarification: these models have two parts

The heart of the relational model can be understood as the combination of the relational data structure with attributes for data values.

Similarly the tree model, as it is typically implemented, combines the tree data structure with descriptive markup node labels (or in XML terminology “generic identifiers”) such as “stanza”.

We will now refer to the tree model as the “tree/DM” model.
The phrase “relational model” does not need the clarification as the role of attributes is commonly understood.

Drum roll

Can we say in general terms why the tree/DM model works so well?

Why it succeeds so well in meeting data management challenges?

We can.

It works the same way the relational model works: *Abstraction* and *Indirection*

Abstraction

Both models focus on the data *itself*, separate from storage and processing. This explicit identification of *data attributes* in one case, and *logical text objects* in the other, brings enormous new functionality and efficiency

Once again: *Indirection*

Both models support an indirect relationship to storage and processing, but in practice the emphasis is often different:

For the relational model abstracting away from storage is dominant

For the tree/DM model abstracting away from processing is dominant

In both case the separation is mediated by a mapping:

logical schema to physical schema in the case of the relational model

text component (type) to processing instructions in the case of trees/DM

Formal vs colloquial understanding of these data models

Although in the relational model we commonly think of attributes as representing dyadic properties or relationships in the world, technically they are names for domains of values,

Similarly in the tree/DM model we think of these node labels as indicating the kind of enclosed text object (stanza, formula, etc), but in the model they are really simply names.

It is our colloquial understanding of these models that enables us to actually use them to secure useful abstraction and data management..

This may seem a small point now, but it motivates a further advance in abstraction, as we well see when we discuss ontologies.

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

**Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.**

Comments and corrections to: renear@illinois.edu.

⑥

IMPLEMENTING THE SOLUTION: XML

Implementing the solution: XML

- XML: a schema language (or: a meta-grammar)
 - XML schemas
 - XML documents
- Example of an XML schema (DTD)
- Example of XML processing
- XML tools
- XML languages

XML

An XML document uses a defined set of delimiters with arbitrary element names and attribute value pairs to nest spans of text.

A *well-formed* XML document fits a formal grammar (along with other constraints) that ensures the document can be parsed as a tree by an XML parser.

NB: a well-formed XML document need not have a schema that defines the element vocabulary and grammar;

it may use arbitrary element, attribute, and value names and arrange text objects in any way that does not violate the tree data structure.

The two main things in the XML world

Schemas [such as Document Type Definitions (DTDs)]

- One for each document type (class, category, genre)
- Defines a markup language for document structures by specifying its vocabulary and syntax (grammar)
- What elements can occur in documents of a particular type, what patterns these elements may form, what other information can be included about these elements?

Document Instances

- Particular documents, marked up with a markup language that meets well-formedness constraints, and, perhaps, also meets the constraints of a relevant schema.

A well-formed XML document

```
<anthology>
  <poem>
    <heading>THE SICK ROSE</heading>
    <stanza>
      <line>O Rose thou art sick.</line>
      <line>The invisible worm,</line>
      <line>That flies in the night</line>
      <line>In the howling storm:</line>
    </stanza>
    <stanza>
      <line>Has found out thy bed</line>
      <line>Of crimson joy:</line>
      <line>And his dark secret love</line>
      <line>Does thy life destroy.</line>
    </stanza>
  </poem>
  <!-- more poems go here -->
</anthology>
```

**Example from the Text Encoding Initiative (P5)*



Schemas for Trees

This is an XML Document Type Definition (DTD), defining an XML document type

```
<!ELEMENT anthology (poem+)>
<!ELEMENT poem (title?, stanza+)>
<!ELEMENT title (#PCDATA)
<!ELEMENT stanza (line+)>
<!ELEMENT line (#PCDATA)>
```

This schema specifies element vocabulary and grammar

The DTD schema language is based on (extended) Backus Naur Form (BNF) grammars

Some XML schema languages provide additional validation and constraints on content, including data typing.

Another DTD

```
<!ELEMENT poem      (title, author? verse) >  
  
<!ATTLIST poem  
          editor      CDATA          #REQUIRED>  
  
<!ELEMENT verse     (stanza+)>  
  
<!ELEMENT stanza    (line+)>  
  
<!ELEMENT title     (#PCDATA | italic | persname)*>  
  
<!ELEMENT author    (#PCDATA)  
  
<!ATTLIST author  
          sex        (male | female)   #IMPLIED  
          dates      CDATA          #IMPLIED  
          bio        IDREF         #IMPLIED>  
  
<!ELEMENT line      (#PCDATA | italic | persname)  
  
<!ATTLIST line  
          lang       CDATA          "ENGLISH">  
  
<!ELEMENT italic    (#PCDATA)>  
  
<!ELEMENT persname (#PCDATA)>
```

Another XML document, with attribute/value pairs

```
<!DOCTYPE text SYSTEM "poem.dtd">

<poem editor="Sara Porter">

<title>Terence</title>
<author person="N320">A. E. Houseman</author>

<verse>
<stanza>
<line>Terence this is stupid stuff </line>
<line>you eat your victuals fast enough</line>
<line>there can't be much amiss 'tis clear</line>

</stanza>
<stanza>
[...]
<line lang="latin">The old lie: </line>
<line> in vino veritas </line>
</stanza>
</verse>

</poem>
```

Valid XML Documents

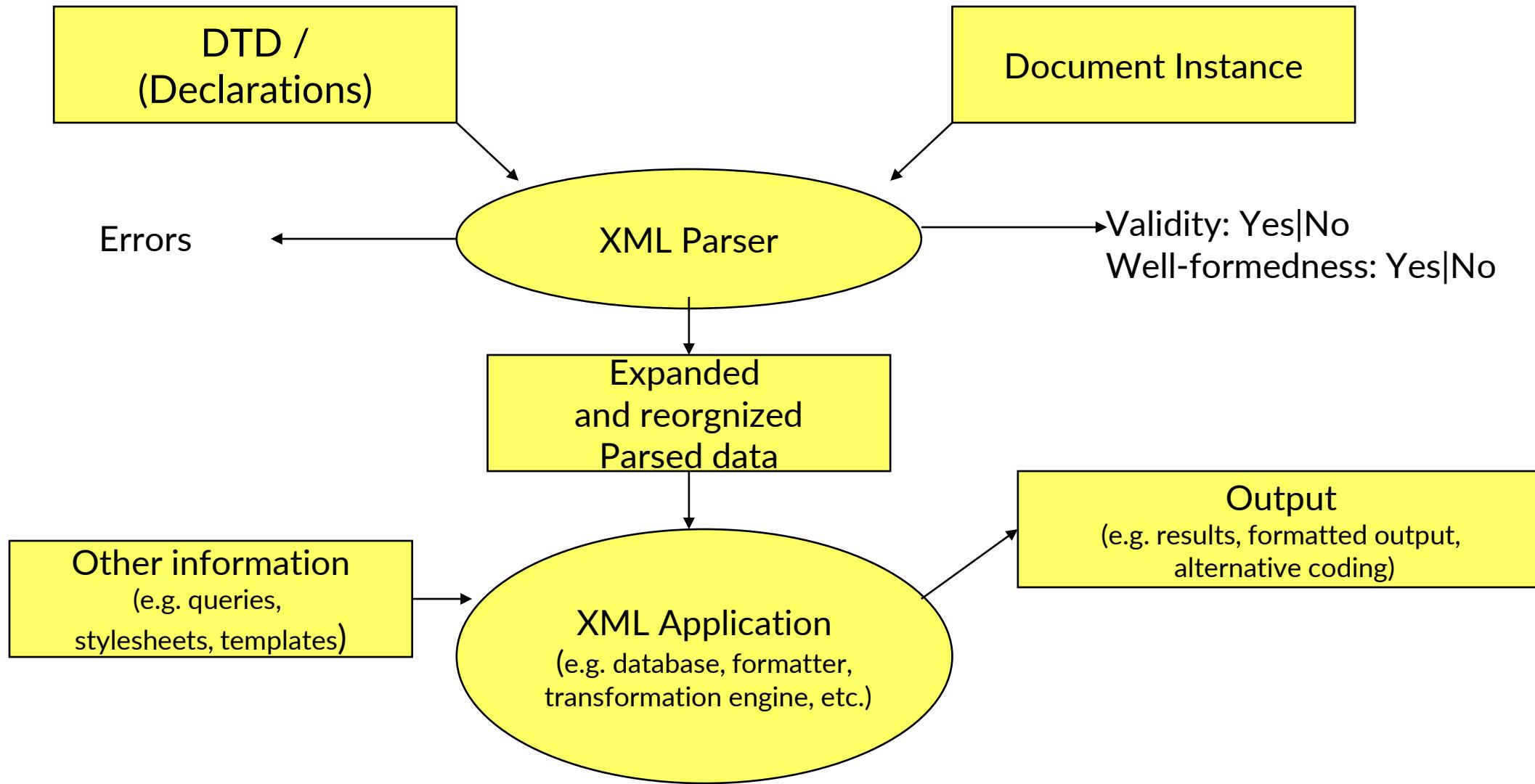
A document instance is *valid* with respect to some schema if it conforms to the declarations in that schema, which is to say, matches the grammar and other constraints.

... nothing out of place, nothing missing, no attributes with values the wrong type, no references that fail, and so on.

A *validating parser* applies an XML schema to an XML document and determines whether or not the document conforms to the constraints specified in the schema.

All valid XML documents are well-formed, but not vice versa.

XML Processing



Some XML tools and schema languages

Two important XML transformation tools

XSLT: “a language for transforming XML documents into other XML documents”

<https://www.w3.org/TR/xslt>

Xquery: “a standardized language for combining documents, databases, Web pages, and almost anything else” <https://www.w3.org/XML/Query/>

Two other XML schema languages

XML Schema (XSD): “the XML Schema Definition Language...offers facilities for describing the structure and constraining the contents of XML documents”

<http://www.w3.org/XML/Schema>

A more complex schema language than DTDs, but does more than validate

Written in XML

Common for business applications

RelaxNG: “a schema language for XML” <http://relaxng.org/>

Similar expressiveness to XSD, with simpler syntax

Less commercial application support



Some important XML languages for documents

You should be familiar with these. Please explore the websites.

XHTML: “a family of current and future document types and modules that reproduce, subset, and extend HTML”

<https://www.w3.org/TR/xhtml1/>

TEI: Text Encoding Initiative – “a standard for the representation of texts in digital form”

<http://www.tei-c.org/index.xml>

JATS: Journal Article Tag Suite “defines a set of XML elements and attributes for tagging journal articles”

<https://jats.nlm.nih.gov/>

XML Languages and Interchange

There are an enormous number of XML markup languages.

See https://en.wikipedia.org/wiki/List_of_XML_markup_languages

Most of these languages are not for text, but are interchange and preservation formats for structured data.

XML is an important preservation format. It use of simple ASCII text with inline tags, can be parsed without a schema, and if a schema is available can be validated ensuring a correct grammar (nothing missing, nothing out of place) and data typing.

FOUNDATIONS OF DATA CURATION (IS531)

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales

School of Information Sciences

University of Illinois at Urbana-Champaign

Includes material adapted from work by Carole Palmer, Melissa Cragin,
David Dubin, Karen Wickett, Bertram Ludaescher, Ruth Duerr and Simone Sacchi.

Comments and corrections to: renear@illinois.edu.