

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

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:

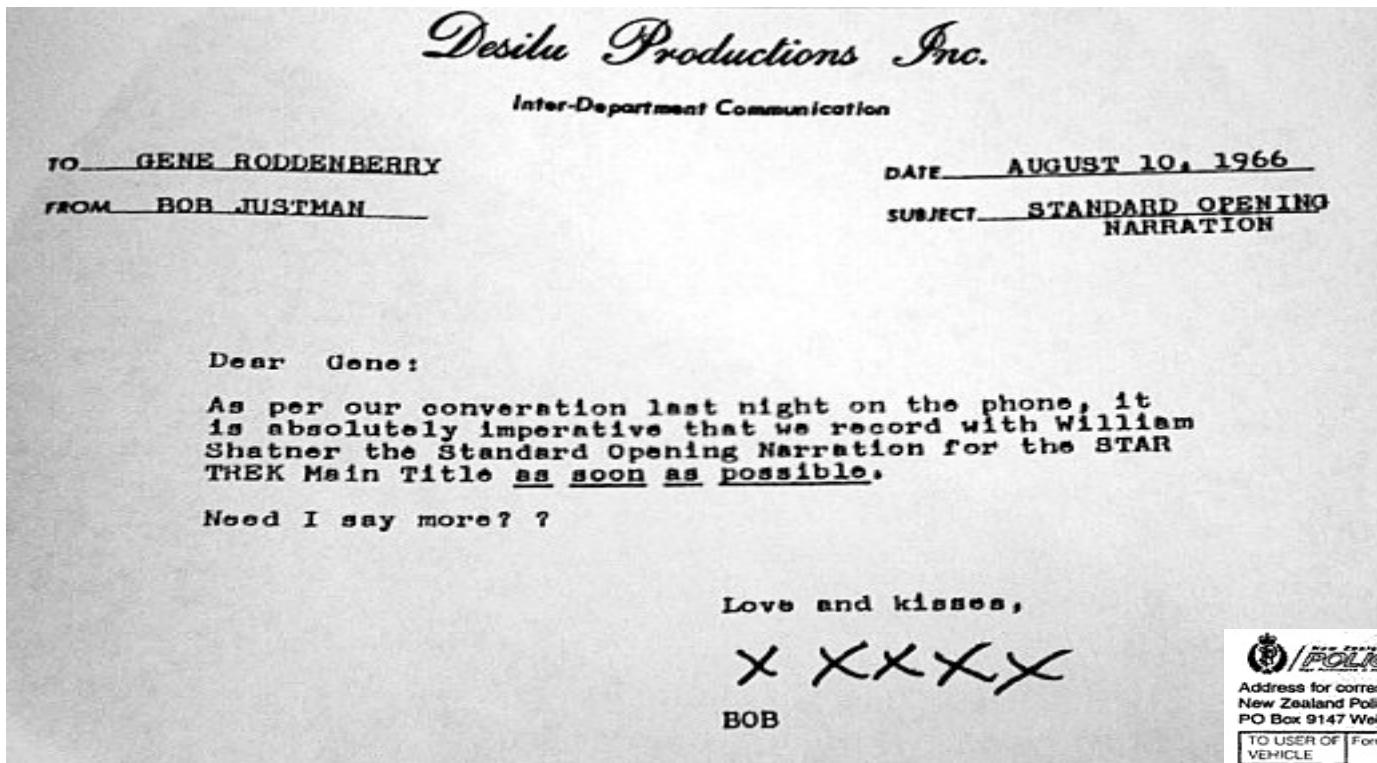
Particularly where the data has this form: *something* has a *value* for an *attribute*

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

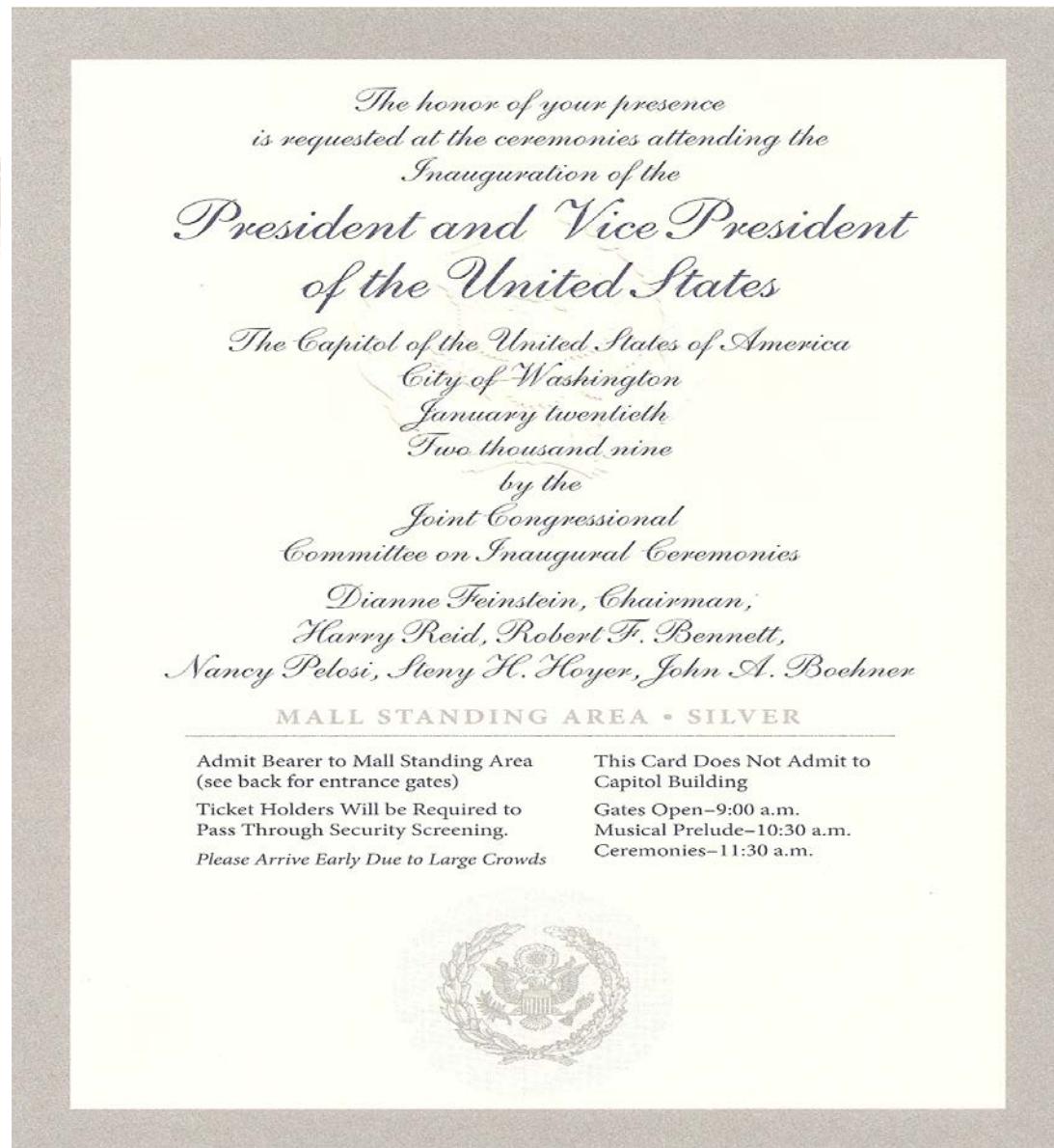
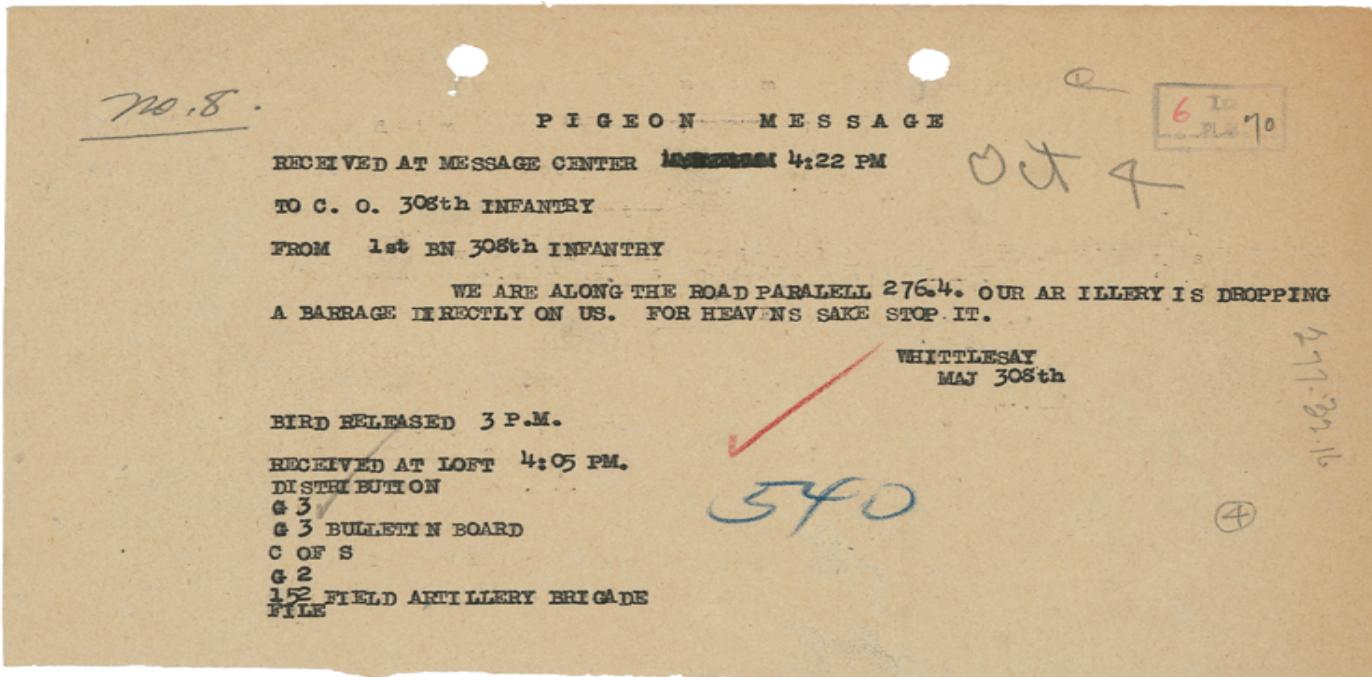
INFRINGEMENT NOTICE
(ISSUED UNDER THE AUTHORITY OF THE LAND TRANSPORT ACT 1968)

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 <i>\$ 120 -</i>
Offence Number <i>1</i>	Offence <i>Exceeded 100 mph</i>	

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 _____

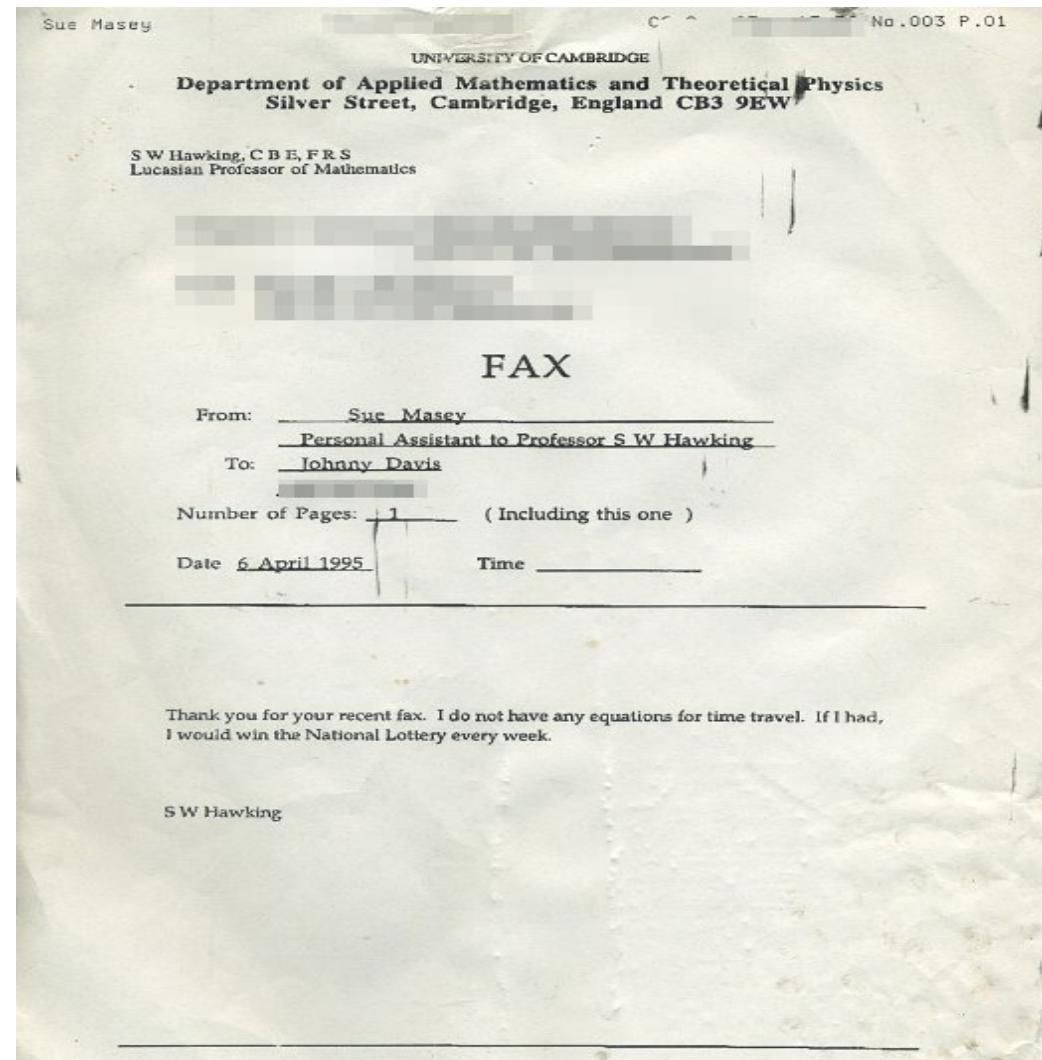
This certificate is issued in pursuance of section 65 of the Marriage Act 1949. Sub-section 5 of that section provides that any certified copy of an entry purporting to be sealed or stamped with the seal of the General Register Office shall be received as evidence of the marriage to which it relates 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 $\langle C_i C_j \rangle$ for a reaction involving two gas phase components i and j . The second-order equation describing the flux term takes the form:

$$\frac{\partial \langle w C_i \rangle}{\partial z} = -w \frac{\partial C_i}{\partial z} - \frac{\partial \langle w C_i \rangle}{\partial z} - \frac{\partial \langle w w C_i \rangle}{\partial z} \quad (2)$$

$$+ \frac{A}{\delta} \cdot \langle C_i \rangle^2 + \frac{1}{\delta} \cdot \langle C_i \rangle^2 + \langle w w C_i \rangle \quad (2)$$

Here term (1) is the diffusional transport term; (2) is the vertical gradient of the turbulent flux; (3) is the change due to suspension; (4) is known as the pressure gradient term and represents correlations of pressure fluctuations with concentration; (5) and (6) terms represent the effect of change in flux due to chemical reactions.

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

$$\frac{\partial \langle w C_i \rangle}{\partial z} = \frac{1}{K} \left(\langle w^2 \rangle - \frac{\partial \langle C_i \rangle}{\partial z} \right) + \langle w R_{C_i} \rangle - \frac{A}{\delta} \langle C_i \rangle^2 - \frac{P_i G_i}{\delta} \quad (2)$$

$$\langle w R_{C_i} \rangle = K_p \langle w C_i \rangle - K_p \langle C_i \rangle + \langle w C_j \rangle - K_p \langle C_j \rangle \quad (2)$$

where the respective correlation terms are defined as in (2). 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:

$$\langle w C_i \rangle = \frac{2k}{A} \langle w^2 \rangle - \frac{\partial \langle C_i \rangle}{\partial z} \quad (2)$$

Here A is a constant derived from rearranging (2). This constant depends on the flow conditions and chemical terms discussed above.

The transport terms in the chemical tendency equations are solved by using a fully implicit finite difference scheme. The vertical grid employed has 45 levels laid out in a logarithmic axis to give highest resolution in the lowest 50 m of the model. The model is solved with a time step of 5 s for the chemistry and 5 s for the dynamics equations. The preliminary focus of these calculations is on evaluating the effect of chemistry on the calculated flux of NO from soil.

Preliminary results

In this initial set of calculations, the mixing coefficient defined as described by (2) were employed in solving the transport equation for NO, NO_x, and CO, keeping 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 defined for the temperature and are referred to as ‘‘no reaction’’ in the following discussions. In both of these calculations, the

Technical Report Documentation Page			
1. Report No.	2. Government Accession No.	3. Recipient's Catalog No.	
FHWA-HEP-12-046			
4. Title and Subject	5. Report Date		
Asset Sustainability Index: A Proposed Measure for Long-Term Performance	July 2012		
7. Author(s)	8. Performing Organization Code		
Gordon D. Proctor, Shobna Varma, Steve Varnedoe	9. Performing Organization Report No.		
9. Performing Organization Name and Address	10. Work Unit No. (TRBID)		
Gordon Proctor & Associates, Inc.	Starzic Corp.	3737 Woodstone Drive	
Dublin, Ohio 43016	Lewis Center, Ohio 43035	7825 Whittle Drive	
National Center for Pavement Preservation			
10. Contract or Grant No.	11. Sponsoring Agency Code		
DTFH61-10-C-00036			
12. Sponsoring Agency Name and Address	13. Type of report and period covered		
FHWA Surface Transportation Environment Planning Cooperative Research Program and the Office of Asset Management, Pavements and Construction	2011		
14. Sponsoring Agency Code			
15. Supplementary Notes			
16. Abstract			
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.			
17. Key Words	18. Distribution Statement		
Asset Sustainability, Asset Management, Long-term Performance, Sustainable Infrastructure Performance Management	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 www.fhwa.dot.gov/hep/step/infrastructure/asstmgmt/		
19. Security Classif.(of this report)	20. Security Classif. (of this page)	21. No. of Pages	22. Price
Unclassified	Unclassified	116	Free

Form DOT F 1700.7 (8-70) Reproduction of completed page authorized



NEUTRAVODNÉ

M9-0144G-50420-00

Slaňovací rám SR-1

Provozní kontroly

Slaňovací rám SR-1 - Provozní kontroly

tabulek

	Strana
1. Odkazy	1
2. Požadovaný stav	1
3. Požadavek na pracovní síly	1
4. Požadovaný stav	4

obrázků

	Strana
1. Kontrola ovládání sláňovacího rámu	2
2. Kontrola ovládání háku TYLER 442 s otevírací pálkou T-Handle Release Unit	3

Odkazy

4. Effectively axiomatized theories

důležitost dokumentace

Požadavky

aný stav

ek na pracovní síly

Skupina Mechanik

ini opatření

Před začátkem

po

M9-0144G-50420-00

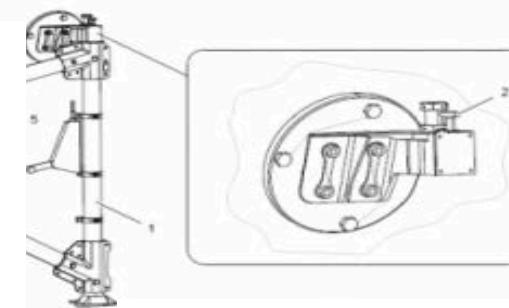
NEUTRAVODNÉ



Postup

1. Provozní kontrola funkčnosti sláňovacího rámu SR-1

- 1.1 Ujistěte se, že SR (Obr. 1 [1]) je v transportní poloze.
- 1.2 Vytahněte držadlo (Obr. 1 [2]) a rukou umístěte SR (Obr. 1 [1]) do sklopné polohy.
- 1.3 Uvolněte držadlo (Obr. 1 [2]):
– Držadlo (Obr. 1 [2]) se musí automaticky zasunout ve sklopné poloze SR.
- 1.4 Ujistěte se, že SR (Obr. 1 [1]) je zajistěn ve sklopné poloze.
- 1.5 Zvedněte konzoli závěsu (Obr. 1 [3]) do horní polohy a zajistěte čep kluzáku (Obr. 1 [4]):
– SR (Obr. 1 [1]) je v pracovní poloze.
- 1.6 Podílné konzole závěsu (Obr. 1 [3]), demontujte čep kluzáku (Obr. 1 [4]) a umístěte konzolu závěsu (Obr. 1 [3]) do spodní polohy a čep kluzáku (Obr. 1 [4]) do držáku (Obr. 1 [5]):
– SR (Obr. 1 [1]) je ve sklopné poloze.
- 1.7 Vytahněte držadlo (Obr. 1 [2]) a rukou umístěte SR (Obr. 1 [1]) do transportní polohy.
- 1.8 Uvolněte držadlo (Obr. 1 [2]):
– Držadlo (Obr. 1 [2]) se musí automaticky zasunout v transportní poloze SR.
- 1.9 Ujistěte se, že SR (Obr. 1 [1]) je zajistěn v transportní poloze.



568
NOČS
(NTI): things identical at one time cannot be distinct at some other time. Sharvy offers a proof that NTI follows from IL (1968, p. 311).

Armed with these principles, Sharvy argues as follows. C-SC and C-69SC had the same members in 1969; therefore, by PE, they were identical in that year. So we have:

(1) In 1969, C-SC = C-69SC.
From (1) and NTI we next deduce that C-SC and C-69SC are identical *simpliciter*:

(2) C-SC = C-69SC.
Now the assumption about C-69SC was:
(3) C-69SC never changes in membership,
and this in combination with (2) yields
(4) C-SC never changes in membership.

That is, the supposedly variable class cannot vary after all. Q. E. D.

This is fine as far as it goes, but we are left to wonder why premise (3) is true. No doubt it is true, but *why?* Isn't just the question? If we are asked to prove that something is true, we must show that it is true. But what does it mean for all classes to have to be invariant if some are, but why must classes such as C-69SC be invariant in the first place?

A similar problem would arise if we tried to argue that Sharvy's strategy to answer my question. The assumption would presumably use in place of NTI the principle of *no contingent identities* (NCI): things identical in one possible world cannot be distinct in another. And in place of (3) it would use the assumption that there are certain worlds that share the same membership in all possible worlds. But what is the ground for this assumption? That is the very thing we want explained.

Whatever the answer, it will have to involve principles in addition to PE and NCI. To see this, consider another PIs, the *Identity of Indiscernibles*:

(ID) *Ind(x,y) (F(x) → y has F) → x = y*.

This is analogous to PE, so if PE and NCI implied that a set has all its members essentially, ID and NCI ought likewise to imply that a set has all its properties essentially. 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.