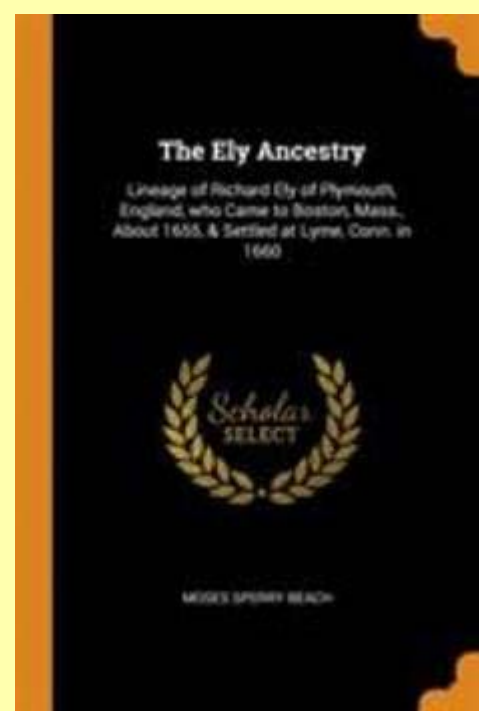
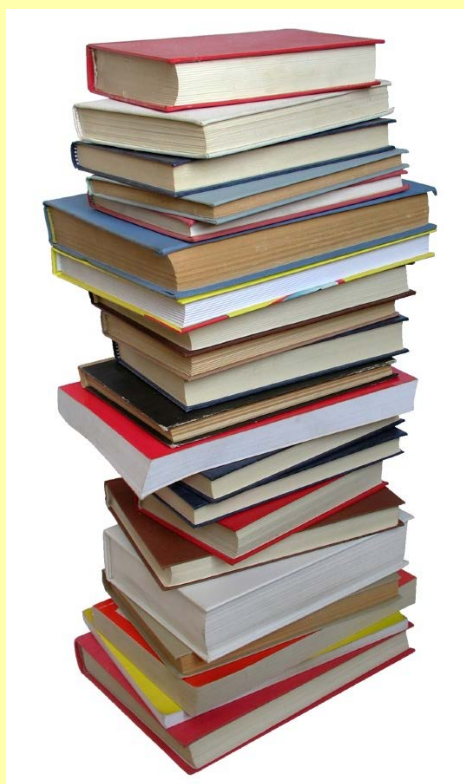
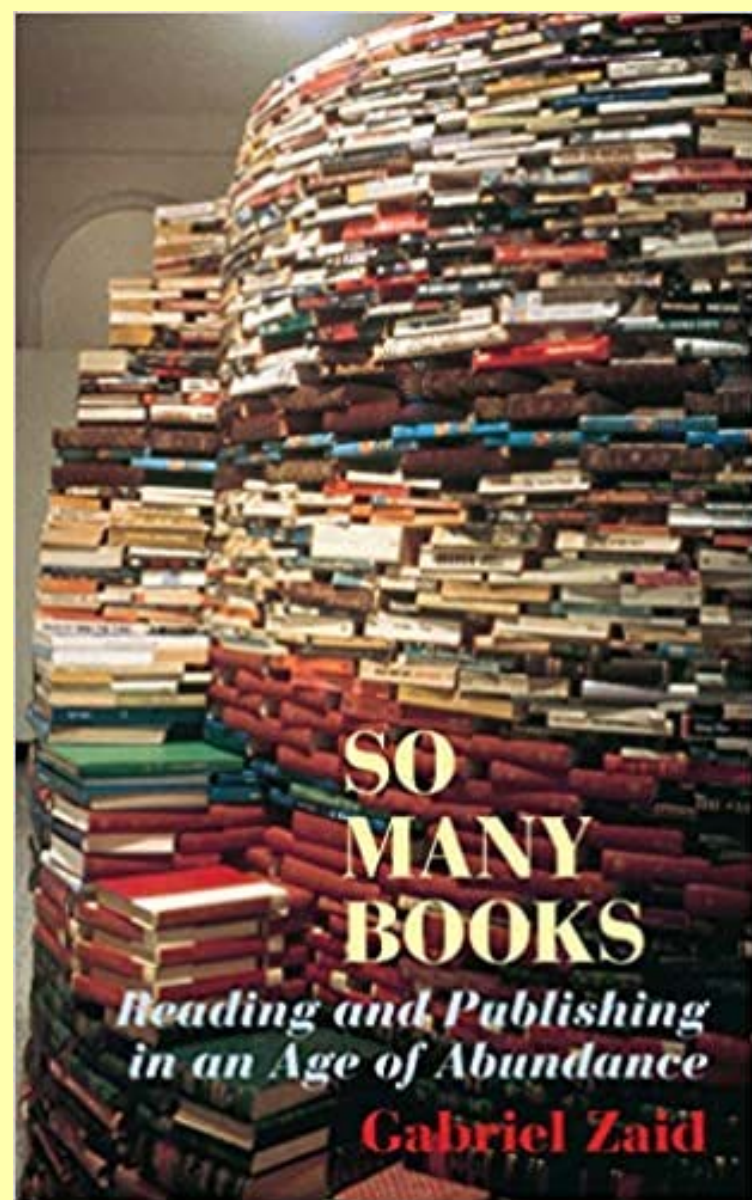


# Near-perfect Relation Extraction from Family Books

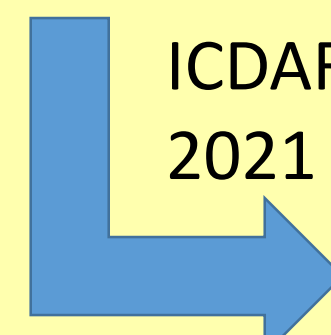
George Nagy

Rensselaer Polytechnic Institute, Troy, NY 12180, USA



GreenEx (NER)

...HEAD:;578,10,9,3,Henrietta,Mills,Hyde,B\_DATE:;578,11,9,1,1826, PARENT1:;578,11,23,2, Julia,Ely,PARENT2:;578,11,37,2,Zabdal,Hyde, M\_DATE:;578,11,54,1, 1848,SPOUSE:;578,11,60,3, Charles,Smith,Shelton,B\_DATE:;578,12,47,1,1819, D\_DATE:;578,12,56,1, 1879,PARENT1:;578,13,1,2, George,Shelton,PARENT2:;578,13,20,2,Betsy,Wooster, CHILD:;578,14,4,2,Fany,Arabella,B\_DATE:;578,14,23,1, 1850, M\_DATE:;578,14,32,1,1874, SPOUSE:;578,14,38,3,Arthur,Harry,Bissell, CHILD:;578,15,4,2, Julia,Elizabeth,B\_DATE:;578,15,24,1,1851, M\_DATE:;578,15,34,1,1878, SPOUSE:;578,15,40,5,Chas.,J.,Van,Tsel, CHILD:;578,16,4,2,Charles,Henry,B\_DATE:;578,16,22,1,1854, CHILD:;578,17,4,2Henry,Hyde,B\_DATE:;578,17,19,1, 1858,HEAD:;



ICDAR 2021 <HEAD [SPOUSE] 57810 9, 578 11 60 Henrietta,Mills,Hyde [Charles,Smith,Shelton]>  
<CHILD [SPOUSE] 578 15 4, 578 15 40 Julia,Elizabeth [Chas., J.,Van Tassel]>

- A family record (above) and two binary relations (below)

## Results:

### Semi-structured Text (Definition)

- Every sign and its must be in the same family record
- A value cannot precede its sign
- Signs and values alternate, except for *collocations* where the sign and value share some or all tokens
- Relations of the same type cannot be nested

Table 1. Data characteristics and accuracy

	Kilbarchan	Miller	Ely
Pages processed	139	389	301
Specified labels	8	10	8
Tokens assigned specified labels	39203	91633	39440
Tokens labeled NONE	34077	131462	101485
Test Set with ground truth			
Pages	6	6	6
Tokens (including "NONE")	3126	3842	3423
Precision	0.999	0.999	0.997
Recall	0.981	0.991	0.992
F-measure	0.990	0.996	0.994

Table 2. Summary of Relation Extraction

	Kilbachan	Miller	Ely
Types to extract	26	26	26
Types found	15	17	16
Families	2615	4186	1219
Extracted Groups	14152	38146	20781
Extracted Relations	11487	25633	18780

### Terminology

Sign = Marker, Query , Search phrase

Value = Extract, Target

Template = (Sign, Value)

### Pseudocode for Named Relation Extraction (NRE)

**Function** Grelex (FamilyRecords, DesiredTuples)

*Input:* FamilyRecords, DesiredTuples

*Output:* ExtractedRelations

Convert FamilyRecords to Families % Family Records is a book-length list of labeled groups of tokens

**For** Family in Families: % Family is a list of Groups, each [Label, Page, Line, Offset, Value], in a single Family; Families is a book-length list of Family(s)

**For** XGroup in Family % restrict search to this Family

**Excerpt** XLabel, XValue, XID **from** XGroup

% XID\_ is Page, Line, Offset of XGroup's value

**For** Tuple in DesiredTuples: % for every specified relation

**Excerpt** LeftLabel, RightLabel **from** Tuple % set the two search arguments

StopGroup ← Stopper(LeftLabel, RightLabel)

% StopGroup is a context-dependent label

**If** XLabel =LeftLabel % if an entity label matches left label of this tuple

**For** YGroup in Family **from** XGroup **to** StopGroup :

% search forward to find a match for the right label

**If** YLabel = RightLabel % if a match for right label is found

**Excerpt** YLabel, YValue, YID **from** YGroup

Relation ← [[XLabel, YLabel, XID YID],[XValue],[YValue]]

% extract this instance of the desired relation

**Append** Relation to ExtractedRelations

**End If** YLabel = RightLabel

**End For** YGroup

**End If** XLabel =LeftLabel

**End For** Tuple

**End For** XGroup

**End For** Family

**Return**(ExtractedRelations) % ExtractedRelations is a list of attributed relations of the form Label\_1[Label\_2], , ID\_1, ID\_2, Value\_1[ Value\_2]

Table 3. Number of extracted relations of each type

Relation	Kilbarchan	Miller	Ely
HEAD[CHILD]	4076	837	3448
HEAD[TWINS]	47	0	0
HEAD[SPOUSE]	2102	1730	1306
HEAD[B_DATE]	6	2849	1128
HEAD[M_DATE]	963	0	1020
HEAD[D_DATE]	0	4119	374
HEAD[BU_DATE]	0	3524	0
HEAD[B_PLACE]	0	2798	0
HEAD[M_PLACE]	134	0	0
HEAD[BU_PLACE]	0	4053	0
HEAD[PARENT1]	0	2824	1203
HEAD[PARENT2]	0	2726	1210
SPOUSE[B_DATE]	3	9	1003
SPOUSE[M_DATE]	12	0	136
SPOUSE[D_DATE]	0	8	360
SPOUSE[B_PLACE]	0	10	0
SPOUSE[M_PLACE]	12	0	0
SPOUSE[PARENT1]	0	9	1015
SPOUSE[PARENT2]	0	7	1028
CHILD[B_DATE]	1143	43	3282
CHILD[C_DATE]	2873	0	0
CHILD[D_DATE]	0	82	831
CHILD[SPOUSE]	34	5	819
CHILD[M_DATE]	38	0	617
TWINS[B_DATE]	6	0	0
TWINS[C_DATE]	38	0	0
<b>TOTAL</b>	<b>11487</b>	<b>25633</b>	<b>18780</b>

Part of a family record generated by Generalized Template Matching:

Sign 1 \_\_\_\_ Value 1 \_\_\_\_ Sign 2 \_\_\_\_ Value 2 \_\_\_\_ Sign 3 \_\_\_\_ Value 3

Binary relations extracted by Named Relation Extraction (NRR):

Sign 1 [Sign 2] = Value 1 {Value 2}

Sign 1 [Sign 3] = Value 1 {Value 3}

Sign 2 [Sign 3] = Value 2 {Value 3}

N-ary relations by combining binary relations:

<HEAD [CHILD, [B\_DATE]] 57810 9, 578 15 40, 578 15 24;  
Henrietta,Mills,Hyde [Julia,Elizabeth [1803]]>

<HEAD, [SPOUSE [[PARENT1], [PARENT2]]] 400 13 08, 400 14 32, 400 15  
27, 400 15 48; Henrietta,Mills,Hyde [Charles,Smith,Shelton  
[[George,Shelton], [Betsy,Wooster]]]>

Attributed relations:

In our implementation, each extracted relation includes the unique ID, Page, Line, and Offset of every participating token.

Acknowledgment:

To my long- time friend, colleague and collaborator Emeritus Professor Dave E. Embley of Brigham Young University, and his team