



南京大學  
NANJING UNIVERSITY



# ESSTER at the EYRE 2020 Entity Summarization Task

Qingxia Liu, Gong Cheng, and Yuzhong Qu

State Key Laboratory for Novel Software Technology, Nanjing University, China

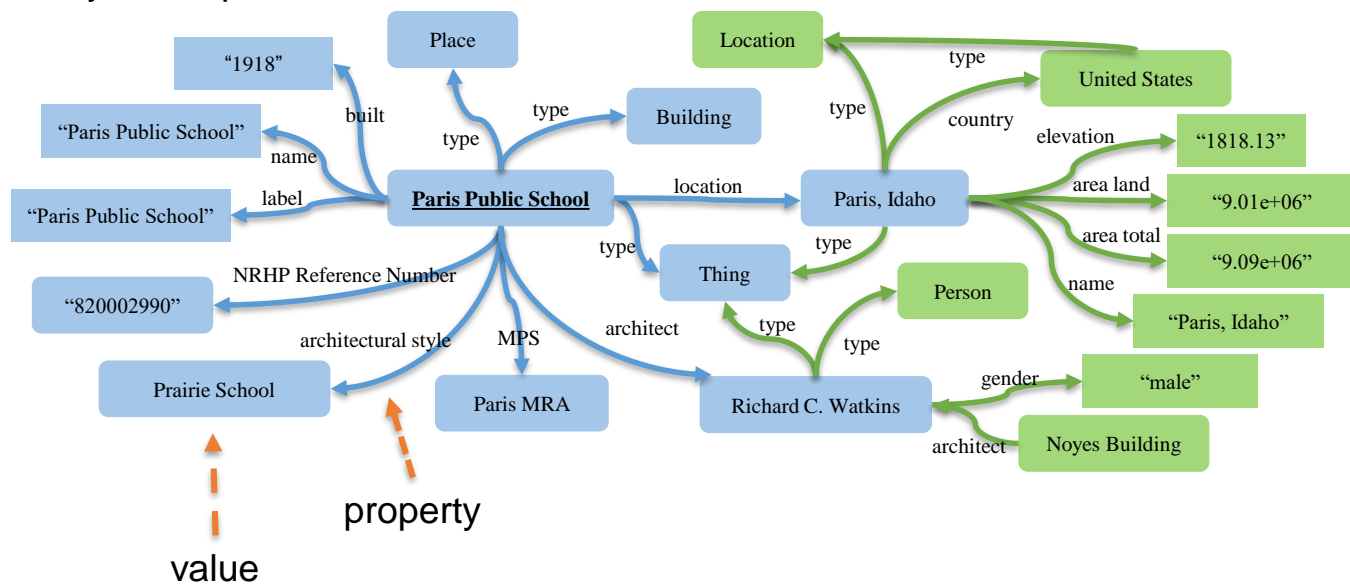
# Outline

- Introduction
- ESSTER
  - Structural Importance
  - Textual Readability
  - Information Redundancy
  - Combinatorial Optimization
- Evaluation
- Conclusion

# Introduction

## RDF Graph

### Entity Description



### Entity Summary

#### **Paris Public School**

type: Building  
location: Paris, Idaho  
built: "1918"  
architect: Richard C. Watkins  
architectural style: Prairie School

# Introduction

- Entity Description  $\text{Desc}(e)$ 
  - $\text{Desc}(e) = \{ \langle e, \text{prop}(t), \text{val}(t) \rangle \}$
  - a set of triples in  $T$ , where  $e$  as subject or object
- Entity Summarization  $\text{Summ}(e)$ 
  - $\text{Summ}(e) \subseteq \text{Desc}(e)$  ,  $|\text{Summ}(e)| \leq k$
  - Provide key information & compact
  - For human users  $\Rightarrow$  Reading Experience

## Paris Public School

MPS: Paris MRA

locmapin: "Idaho"

NRHP Reference Number: "82000290"

coord format: "dms"

latitude: "42.2284"

## Paris Public School

type: Thing

type: Place

type: Building

label: "Paris Public School"

name: "Paris Public School"

## Paris Public School

type: Building

location: Paris, Idaho

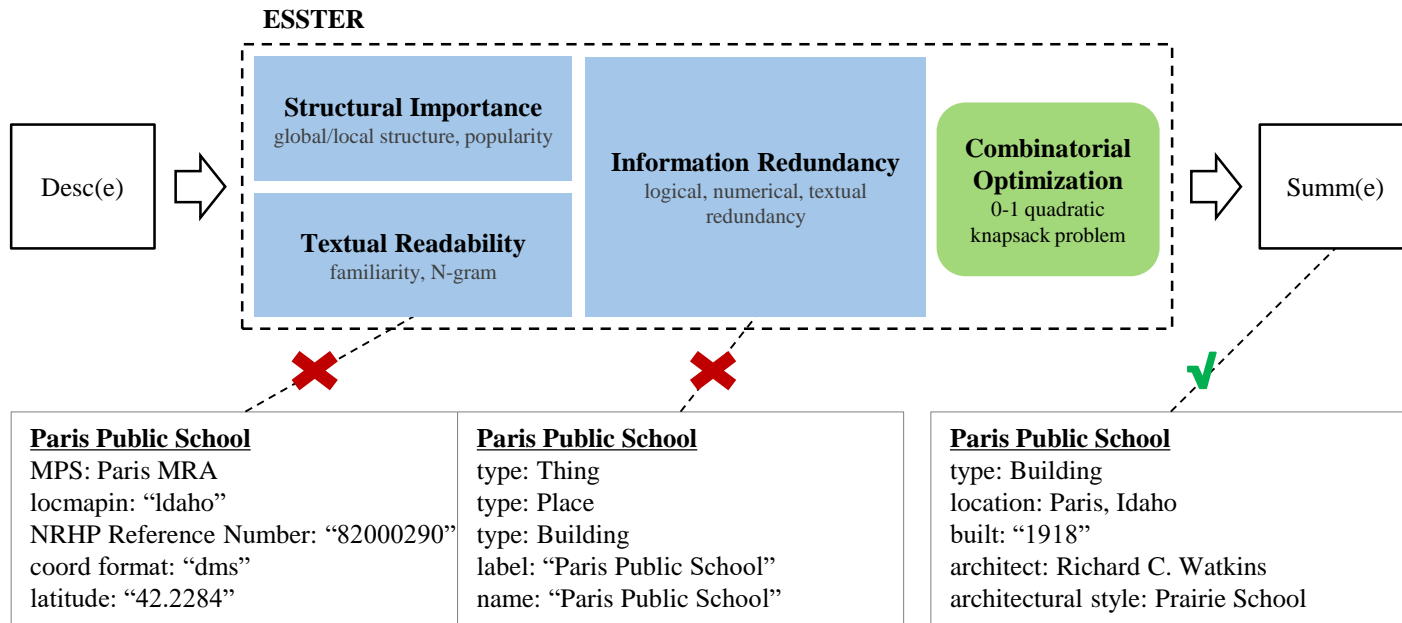
built: "1918"

architect: Richard C. Watkins

architectural style: Prairie School

# ESSTER

Generating entity summaries of structural importance, high readability, and low redundancy.



# Structural Importance

## ■ Structural Importance

- $W_{\text{struct}}(t) = \alpha \cdot \text{glb}(t) + (1 - \alpha) \cdot \text{loc}(t)$

## ■ Global

- $\text{glb}(t) = \text{ppop}_{\text{global}}(t) \cdot (1 - \text{vpop}(t))$
- generality of property
- characteristic of value

type: Thing  
type: Architect ✓  
hypernym: Architect

## ■ Local

- $\text{loc}(t) = (1 - \text{ppop}_{\text{local}}(t)) \cdot \text{vpop}(t)$
- punishment on multi-valued properties
- avoiding too technical/specific values

subject: American architects  
subject: NRHP architects  
birth year: “1858” ✓  
wiki page ID: “34981613”

# Textual Readability

location ✓ NRHP Reference Number ?

birth date ✓  coord format ?

country ✓ locmapin ?

architectural style ✓ MPS ?

# Textual Readability

location ✓ NRHP Reference Number ?

birth date ✓ coord format ?

country ✓ locmapin ?

architectural style ✓ MPS ?



A property is familiar to users if  
it is often used in an open-domain corpus



# Textual Readability

location ✓ NRHP Reference Number ?

birth date ✓ coord format ?

country ✓ locmapin ?

architectural style ✓ MPS ?



$B$  : # documents in the corpus

$b(t)$  : # documents where prop(t) appears

$M$  : # documents have been read by the user

$m$  : # documents have been read by the user  
and prop(t) appears

## ■ Familiarity

$$Q(t) = \sum_{m=0}^{\min(b(t), M)} \frac{\binom{b(t)}{m} \cdot \binom{B-b(t)}{M-m}}{\binom{B}{M}} \cdot \text{familiarity}(m),$$

$$\text{familiarity}(m) = \frac{\log(m+1)}{\log(B+1)}.$$

## ■ Textual Readability

$$W_{\text{text}}(t) = \log(Q(t) + 1).$$

# Information Redundancy

- Redundancy
  - $\text{sim}(t_i, t_j)$
- Logical Redundancy
- Numerical Redundancy
- Textual Redundancy

# Information Redundancy

## ■ Redundancy

- $\text{sim}(t_i, t_j)$

## ■ Logical Redundancy

$\text{sim}(t_i, t_j) = 1$  if

- $\text{prop}(t_i) = \text{prop}(t_j) = \text{rdf:type}$  ,  $\text{val}(t_i)$  and  $\text{val}(t_j)$  have  $\text{rdfs:subClassOf}$  relation
- $\text{val}(t_i) = \text{val}(t_j)$  ,  $\text{prop}(t_i)$  and  $\text{prop}(t_j)$  have  $\text{rdfs:subPropertyOf}$  relation

type: Thing  
type: Place  
type: Building

label: "Paris Public School"  
name: "Paris Public School"

# Information Redundancy

## ■ Redundancy

- $\text{sim}(t_i, t_j)$

## ■ Numerical Redundancy

$$\text{sim}(t_i, t_j) = \max\{\text{sim}_p(t_i, t_j), \text{sim}_v(t_i, t_j), 0\}$$

- $\text{val}(t_i)$  and  $\text{val}(t_j)$  are both numerical values

- $\text{sim}_p(t_i, t_j) = \text{ISub}(\text{prop}(t_i), \text{prop}(t_j))$

- $\text{sim}_v(t_i, t_j) = \begin{cases} -1, & \text{if } \text{val}(t_i) \cdot \text{val}(t_j) \leq 0 \\ \frac{\min\{\text{val}(t_i), \text{val}(t_j)\}}{\max\{\text{val}(t_i), \text{val}(t_j)\}}, & \text{otherwise} \end{cases}$

area land: "9.01e+06"	} <b>sim=0.9019</b>
area total: "9.09e+06"	
elevation: "1818.13"	} <b>sim=0.0002</b>

# Information Redundancy

## ■ Redundancy

- $\text{sim}(t_i, t_j)$

## ■ Textual Redundancy

$$\text{sim}(t_i, t_j) = \max\{\text{sim}_p(t_i, t_j), \text{sim}_v(t_i, t_j), 0\}$$

- $\text{sim}_p(t_i, t_j) = \text{ISub}(\text{prop}(t_i), \text{prop}(t_j))$
- $\text{sim}_v(t_i, t_j) = \text{ISub}(\text{val}(t_i), \text{val}(t_j))$

given name: "Richard"	}	<b>sim=0.78</b>
name: "Richard C. Watkins"		
gender: "male"	}	<b>sim=0</b>

# Combinatorial Optimization

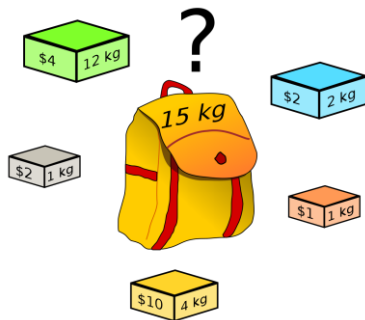
- Goal: generate entity summary of
  - high structural importance  $W_{\text{struct}}(t_i)$
  - high textual readability  $W_{\text{text}}(t_i)$
  - low information redundancy  $-\text{sim}(t_i, t_j)$
  - and satisfy size constraint  $k$

# Combinatorial Optimization

## ■ 0-1 Quadratic Knapsack Problem (QKP)

- max profit, satisfy weight constraint
- profit:
  - structural importance, textual readability
  - low information redundancy
- weight:  $x_i$ 
  - 1 if  $t_i \in \text{Summ}(e)$
  - 0 otherwise

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^{|\text{desc}(e)|} \sum_{j=i}^{|\text{desc}(e)|} \text{profit}_{i,j} \cdot x_i \cdot x_j, \\ & \text{subject to} && \sum_{i=1}^{|\text{desc}(e)|} x_i \leq k, \\ & && x_i \in \{0, 1\} \text{ for all } i = 1 \dots |\text{desc}(e)|. \end{aligned}$$



$$\text{profit}_{i,j} = \begin{cases} (1 - \delta) \cdot (W_{\text{struct}}(t_i) + W_{\text{text}}(t_i)) & i = j, \\ \delta \cdot (-\text{sim}(t_i, t_j)) & i \neq j, \end{cases}$$

# Evaluation

- Data
  - ESBM v1.2
- Baselines
  - 9 unsupervised methods
- Results
  - top-2 on DBpedia under  $k=5$
  - best in all the other three settings

F1 Scores

	DBpedia		LinkedMDB	
	$k = 5$	$k = 10$	$k = 5$	$k = 10$
RELIN	0.242	0.455	0.203	0.258
DIVERSUM	0.249	0.507	0.207	0.358
FACES	0.270	0.428	0.169	0.263
FACES-E	0.280	0.488	0.313	0.393
CD	0.283	0.513	0.217	0.331
LinkSUM	0.287	0.486	0.140	0.279
BAFREC	<b>0.335</b>	0.503	0.360	0.402
KAFCa	0.314	0.509	0.244	0.397
MPSUM	0.314	0.512	0.272	0.423
ESSTER	0.324	<b>0.521</b>	<b>0.365</b>	<b>0.452</b>



# Conclusion

- ESSTER: generating entity summaries by integrating
  - structural importance,
  - textual readability,
  - and information redundancy
  - via combinatorial optimization
- ESSTER achieves SOTA among unsupervised entity summarizers on ESBM v1.2
- Future Work
  - more powerful measures of readability and redundancy
  - incorporate these features into a neural network model



南京大學  
NANJING UNIVERSITY



# Thank you !

Questions ?