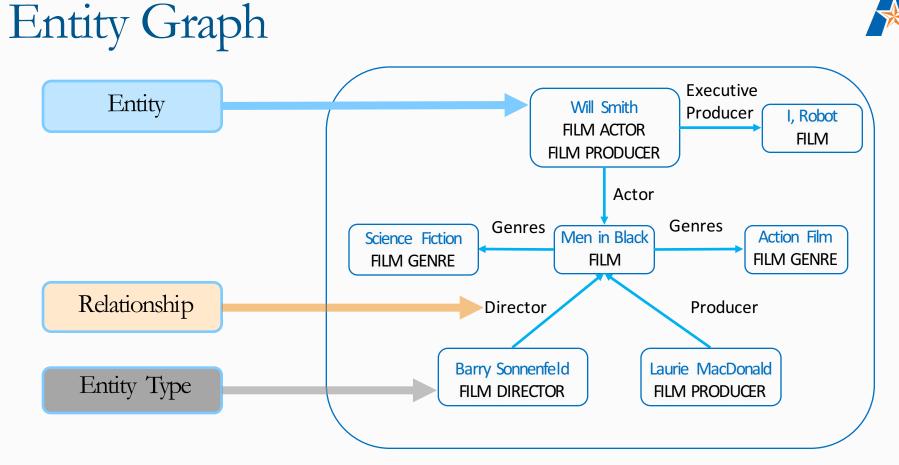
### Generating Preview Tables for Entity Graphs

<sup>#</sup>Ning Yan, \*Sona Hasani, \*Abolfazl Asudeh, \*Chengkai Li
<sup>#</sup>Huawei U.S. R&D Center
\*University of Texas at Arlington, Innovative Database and Information Systems Research

(IDIR) Laboratory

SIGMOD 2016, June 30th, 2016





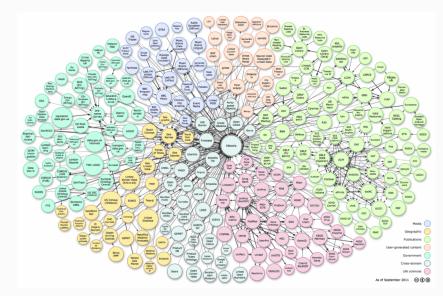


# Ultra-heterogeneous Entity Graphs

Large and complex graphs capturing millions of entities and billions of relationships between entities.

Applications: search, recommendation systems, business intelligence, health informatics, fact checking

Freebase : 1.9 billion triples DBpedia : 3 billion triples YAGO : 120 million triples Linked Open Data : hundreds of datasets 52 billion RDF triples

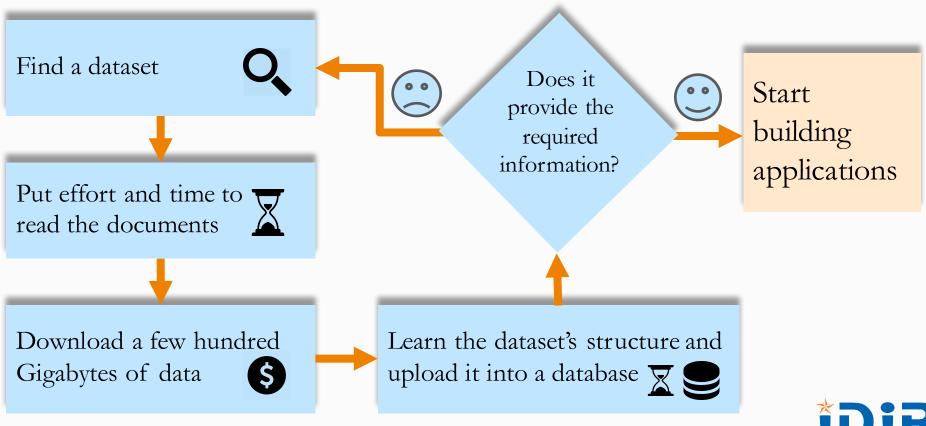


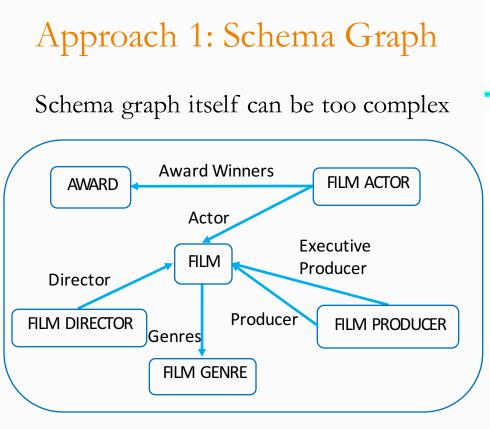
#### http://linkeddata.org/



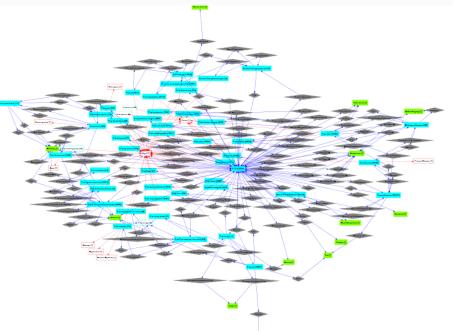
### Steep Flag-Down Cost







Need for a Quick Overview



Schema graph of "Film" domain in Freebase Entity graph: 2M entities, 18 M edges Schema graph: 63 entity types ,136 edges



# Need for Quick Overview



### Approach 2: Schema Summary

Schema summarization in relational database [Yang PVLDB09, Yang PVLDB11]

- Cluster tables in relational database by their semantic roles and similarities.
- o Clusters tables, not relationships
- o Detailed

#### XML summarization [Yu VLDB06]

• Provide a succinct overview of the entire schema graph

#### Graph summarization [Tian SIGMOD08, Zhang ICDE10]

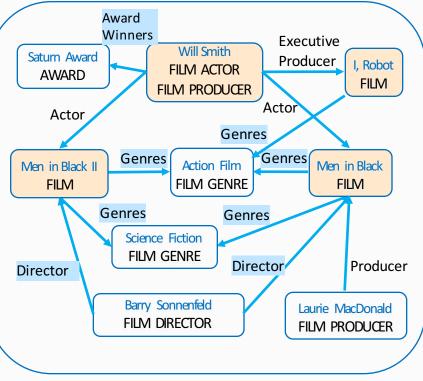
• Group graph nodes based on their attribute similarity and allow users browse the summary from different grouping granularities.



### **Preview Tables**



Key :	attribu	tes		Non	-key attributes	
	FILM A	ACTOR	Award	Winners		
	Will Sm	ith	Saturn A	Award		M
			•			
						Dir
FII	FILM Direc		ctor		Genres	
Men in 1	Black	Barry Sor	nnenfeld	{Action ]	Film, Science Fiction}	
Men in 1	Black II	Barry Sor	nnenfeld	{Action ]	Film, Science Fiction}	
I, Robot	t		-	Action F	ïlm	





# Too Many Previews. Which one to Choose?

- Many possible previews
- 0 Different choices

Preview A	Preview B							
FILM         FILM SET DECORATOR         FILM EDITOR	FILM         FILM DIRECTOR         FILM PRODUCER							
PRODUCTION COMPANY     FILM								





### Score of the Table? Score of the Preview?

FILM	Act	tor	Genres	$4 \times (6+5) = 44$		
4	6	5	5	$1 \land (0 \lor 3) = 11$		
			-			44+16=60
FILM ACTOR	Actor	A	ward Winners			
2	6		2	$-2 \times (6+2) = 16$		



### Key Attributes Scoring

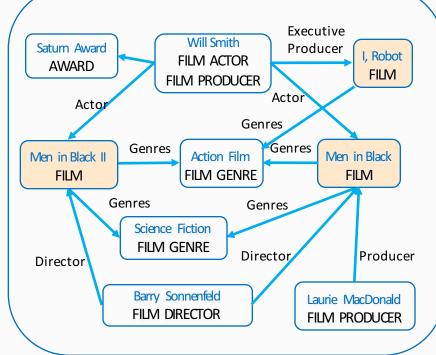


Coverage-based method

Coverage(FILM) = 3

### Random walk-based method

Stationary distribution of a random walk process defined over the schema graph



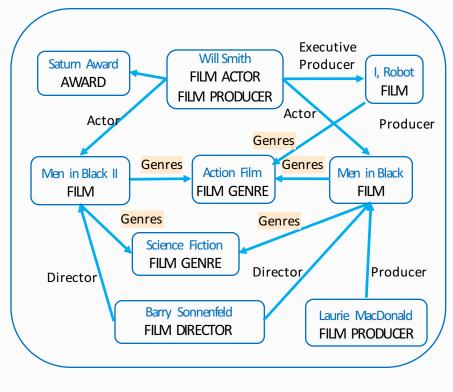


### Non-key Attributes Scoring



Coverage-based method Coverage(Genres) = 5 Entropy-based method Entropy(Genres) =  $(2/3) \log(3/2) + (1/3) \log(3/1)$ = 0.28

FILM	Director	Genres			
Men in Black	Barry Sonnenfeld	{Action Film, Science Fiction}			
Men in Black II	Barry Sonnenfeld	{Action Film, Science Fiction}			
I, Robot	_	{Action Film}			





# Optimal Preview Discovery



Find the preview with highest score that satisfies

- Size constraint
  - Number of key attributes K• Number of non-key attributes N• Distance leaves a terms in tables J• Diverse dist(Ti, Tj)  $\geq d$
- Distance between two preview tables d

FILM	Performance	es Genres	Direc	cted By	]	FILM FESTIVAL	Location	Focus
FILM I	DIRECTOR	Films Direc	ted			FILM COMPANY	Films	
FILM P	RODUCER	ODUCER Films Produced				FILM CHARACTE	R Portra	yed in Film
Tight						Divers	se	



### Concise Preview, Dynamic Programming Algorithm



We assume all K key attributes are ordered arbitrarily.

optimal concise preview (k, n, X) is the best of:

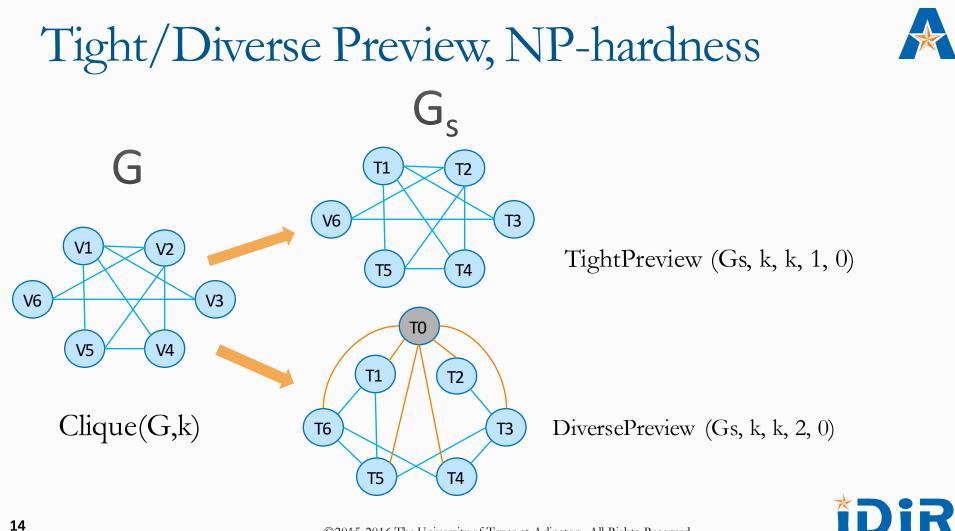
```
optimal concise preview (k, n, X-1)
```

```
optimal concise preview (k-1, n-1, X-1) U X-th Key-attribute with 1 non-key attribute
```

```
optimal concise preview (k-1, n-2, X-1) \cup X-th Key-attribute with 2 non-key attributes ... ...
```

optimal concise preview (k-1, k-1, X-1)  $\cup$  X-th Key-attribute with (n-k+1) non-key attributes





Tight/Diverse Preview, Apriori Property Algorithm



1. Construct 2-cliques by enumerating all key attribute pairs

2. for i = 3 to k

generate i-cliques from (i-1)-cliques based on Apriori property

3. find the k-clique with highest score, **return** as optimal preview







#### Dataset: Freebase

#### Accuracy of scoring measures

- Compared with Freebase ground truth:
  - Key attributes: Precision-at-k (p@k), Average Precision, Discounted Cumulative Gain (nDCG)
  - o Non-key attributes: Mean Reciprocal Rank (MRR)
- Compared with crowd ranking:
  - o Pearson Correlation Coefficient (PCC)

### Efficiency of algorithms

• Execution time

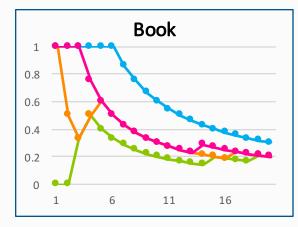
### User study

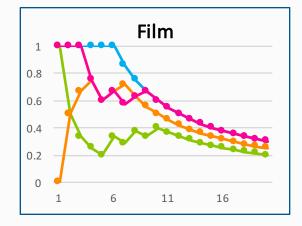
- Existence test questions
- o User experience questions

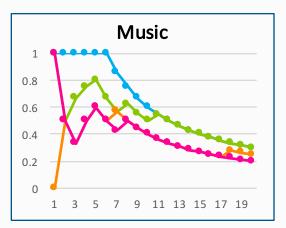


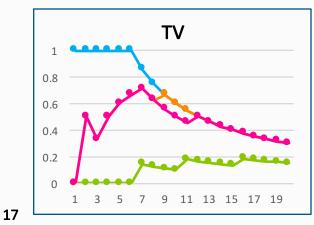
### Key Attribute Scoring (p@k)

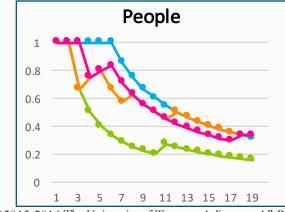


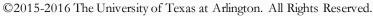












Optimal p@k
YPS09 [Yang PVLDB09]
Coverage
Random Walk

### MRR of Non-key Attributes (MRR)



Domain	Coverage	Entropy
books	0.8	0.786
film	0.2	0.25
music	0.528	0.589
TV	0.622	0.379
people	0.708	0.606



### Compare Rankings by Our Approach and the Crowd (PCC)



		Key Attril	Non-key Attribute		
Domain	YPS09	Coverage	Random Walk	Coverage	Entropy
books	0.4	0.55	0.43	0.43	0.43
film	-0.01	0.48	0.25	0.35	0.35
music	0.37	0.33	0.46	0.42	0.41
TV	0.37	0.69	0.65	0.47	0.47
people	0.36	0.31	0.29	0.43	0.43



#### Domains: film **Key Attributes** Non-Key Attributes 015 Film art director Film crew gig Personal film appearance (45712) Film cut (173046) Film casting director books Crewmember(Film

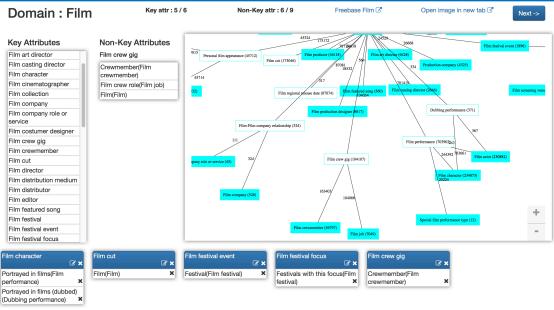
User Study, Hand-crafted Preview Tables

music ΤV people

Ο

- Hand-crafted preview tables 0
- 10 PhD students in Database 0 research group
- Individually and as a group Ο
- \$20 gift card Ο







# User Study, Approaches Compared

- Domains: film, books, music, TV, people
- Approaches:
  - o Schema graph
  - o Concise preview
  - o Tight preview
  - o Diverse preview
  - Freebase ground truth
  - o YPS09
  - Hand-crafted preview tables
- 4 existence questions
- 4 experience questions
- 84 Master's and PhD students in database area
- \$15 gift card

. в	ased on this	schema sum	nmary, I I	know the da	tase	t cont	ains entit	ies that are "fi	lm produc	cer".		
	Agree			Not sure							Next	
0	Disagree											
ilm	View in a new Tab											
+	Film character	Portra actor: Peter McKenzie film: The Lord of the Rin	<b>yed in films</b> gs: The Fellowshi	p of the Ring		Film actor sty Springfie		<b>m performances</b> e: The Best of 1969-1970: Vo		Film directo acobo Moral	r Films directed	
	Foghorn Leghorn	actor: Mel Blanc film: False Hare special_performance_typ	e: Voice		H	Ripken, Jr. Il Smith	film: Baseball the	e Ripken Way: Pitching mmad Ali		Annie Leibov Starhawk	Tull Circle	
	Strawberry Shortcake	actor: Russi Taylor film: Strawberry Shortca special_performance_typ		rykins								
						_						
	Film	Performances	Genres	Country of origin	Runtim	•	Film crewmember	Films crewed	1	- 10		
	ong of the Thin Man	actor: Myrna Loy character: Nora Charles	Thrille	United States of America	runtime		Art Babbitt	Noname	Aamir	producer Khan L	Films Produced	
	The Hudsucker	actor: Tim Robbins		United States of	<u> </u>	1 [	Krishna	Noname	Ralph	Fiennes	oriolanus	
	Proxy	character: Norville Barnes	Comedy	America	runtime		Shelby Young	film: My Soul to Take film_crew_role: Adr Voices	Arthu	r Freed A	n American in Paris	
	l Am Sam	actor: Sean Penn character: Sam Dawson	Courtroom Drama	United States of America	runtime		,		<u>.</u>			

5. Very easy	Additional Comments Next
4. Easy	
3. Neutral	
2. Hard	
1. Very hard	

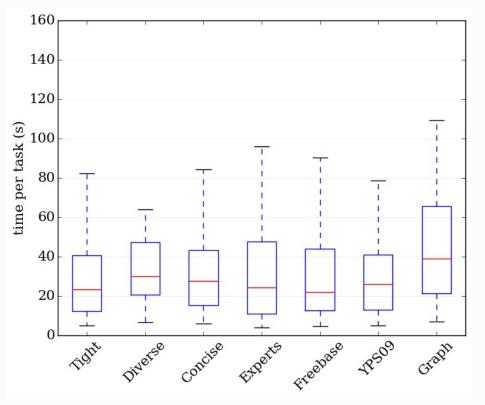
### User Study: Existence Test Questions

	Tight	Diverse	Freebase	Experts	YPS09	Schema Graph
Concise	z=1.59 p=0.0559	z=-2.28 p=0.0113	z=0.49 p=0.3121	z=-0.13 p=0.4483	z=0.36 p=0.3594	z=-0.43 p=0.3336
Tight		z=-3.48 p=0.0003	z=-1.12 p=0.1314	z=-1.69 p=0.0455	z=-1.282 p=0.0999	z=-1.93 p=0.0268
Diverse			z=2.57 p=0.0051	z=2.10 p=0.0179	z=2.60 p=0.0047	z=1.70 p=0.0446
Freebase				z=-0.61 p=0.2709	z=-0.15 p=0.4404	z=-0.87 p=0.1922
Experts					z=0.49 p=0.3121	z=-0.29 p=0.3859
YPS09						z=-0.77 p=0.2206

#### Pairwise comparisons of conversion rates, domain="music", $\alpha$ =0.1



### User Study: Existence Test Questions



Time taken on existence tests, domain="music"



### User Study: User Experience Questions



Questions	most fa	avorabl	e	least favorable			
How easy was it to read the schema summary?	Freebase	Diverse	Graph	Experts	YPS09	Concise	Tight
How much understanding of the data can you gain from it?	Graph	Freebase	YPS09	Diverse	Concise	Tight	Experts
How helpful was it in assisting you to understand the data?	Graph	Freebase	YPS09	Diverse	Experts	Concise	Tight
Is it missing important information?	YPS09	Concise	Experts	Graph	Tight	Freebase	Diverse

Systems sorted by average user experience scores across five domains





# Acknowledgment







**Disclaimer**: This material is based upon work partially supported by the National Science Foundation Grants 1018865, 1408928 and the National Natural Science Foundation of China Grant 61370019. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.



# Thank You! Questions?



### Generating Preview Tables for Entity Graphs

Ning Yan, Sona Hasani, Abolfazl Asudeh, Chengkai Li

