



COLLEGE OF ENGINEERING

COMPUTER SCIENCE & ENGINEERING

UNIVERSITY OF MICHIGAN



MEDICAL SCHOOL

UNIVERSITY OF MICHIGAN



GEMS LAB

Graph Summarization Meets Outlier Detection

Danai Koutra

Morris Wellman Assistant Professor, CSE

Computational Medicine and Bioinformatics (courtesy)

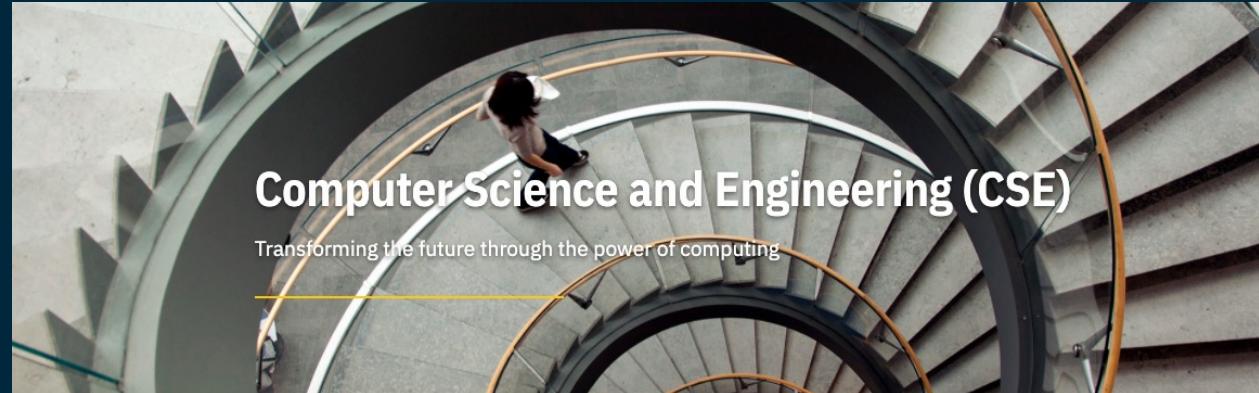
Associate Director for AI, Michigan Institute for Data Science

ACM SIGKDD, 6th Outlier Detection and Description (ODD) Workshop – August 15, 2021

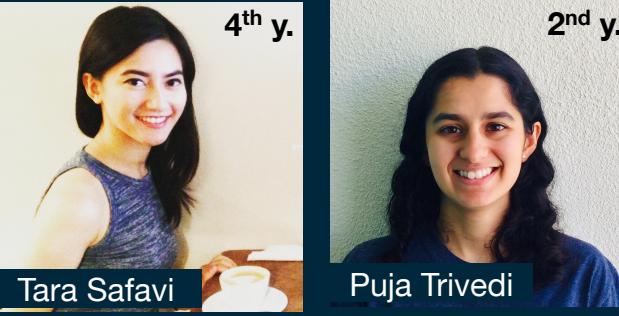
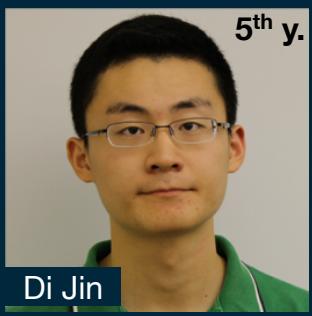
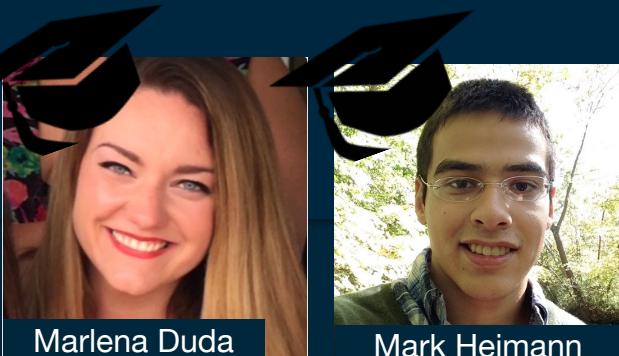
Joint work with: Caleb Belth, Christos Faloutsos, Brian Gallagher, Mark Heimann, Di Jin, Yike Liu, Ryan Rossi, Tara Safavi, Neil Shah, Chandra Sripada, Jilles Vreeken, Xinyi Zheng, ...

About me: Danai Koutra

- Morris Wellman *soon-to-be*
Assoc. Professor in CSE,
University of Michigan (eff. Sep 1)
- Associate Director for AI,
Michigan Institute for Data
Science (MIDAS)



The image shows the homepage of the Michigan Institute for Data Science (MIDAS). The header features the MIDAS logo (a yellow 'M' icon followed by the word 'MIDAS') and the text 'MICHIGAN INSTITUTE FOR DATA SCIENCE UNIVERSITY OF MICHIGAN'. Below the header is a navigation bar with links for HOME, ABOUT, RESEARCH, TRAINING, EVENTS, PARTNERSHIPS, PEOPLE, RESOURCES, and CONTACT US. A search icon is also present. The main content area has a hexagonal grid background and features a section titled 'OUR MISSION' with the text: 'MIDAS strengthens University of Michigan's preeminence in Data Science and Artificial Intelligence, and enables their transformative use in a wide range of research disciplines to achieve lasting societal impact.'



Welcome!

We are the **Graph Exploration and Mining at Scale (GEMS)** lab at the [University of Michigan](#), founded and led by [Danai Koutra](#). Our team researches important data mining and machine learning problems involving interconnected data: in other words, *graphs or networks*.

From airline flights to traffic routing to neuronal interactions in the brain, graphs are ubiquitous in the real world. Their properties and complexities have long been studied in fields ranging from mathematics to the social sciences. However, many pressing problems involving graph data are still open. One well-known problem is *scalability*. With continual advances in data generation and storage capabilities, the size of graph datasets has dramatically increased, making scalable graph methods indispensable. Another is the changing nature of data. Real graphs are almost always *dynamic*, evolving over time. Finally, many important problems in the social and biological sciences involve analyzing not one but *multiple* networks.

So, what do we do?

The problems described above call for **principled, practical, and highly scalable graph mining methods**, both theoretical and application-oriented. As such, our work connects to fields like linear algebra, distributed systems, deep learning, and even neuroscience. Some of our ongoing [projects](#) include:

- Algorithms for [multi-network tasks](#), like matching nodes across networks
- Learning [low-dimensional representations of networks](#) in metric spaces
- Abstracting or “[summarizing](#)” a graph with a smaller network
- Analyzing [network models of the brain](#) derived from fMRI scans
- [Distributed graph methods](#) for iteratively solving linear systems
- Network-theoretical [user modeling](#) for various data science applications

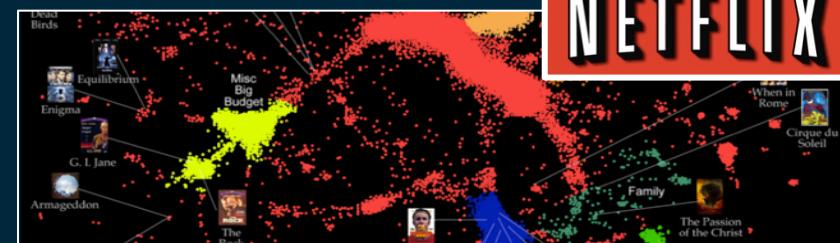
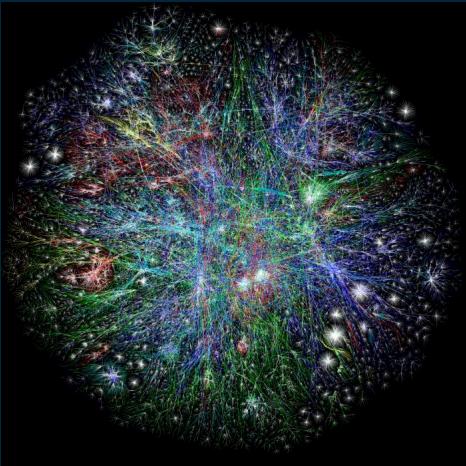
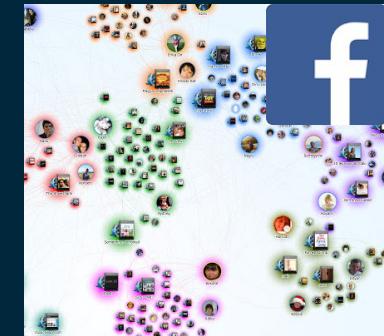
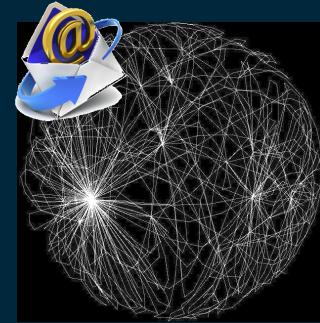
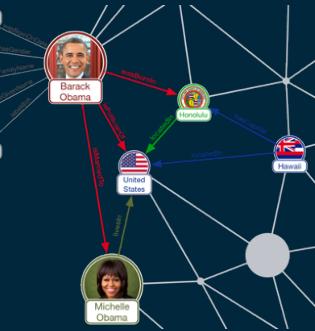
We're grateful for funding from [Adobe](#), [Amazon](#), the Army Research Lab, the Michigan Institute for Data Science (MIDAS), Microsoft Azure, the National Science Foundation (NSF), and

Interested?

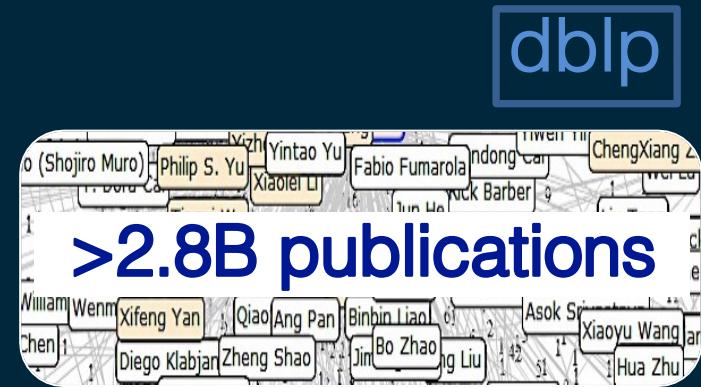
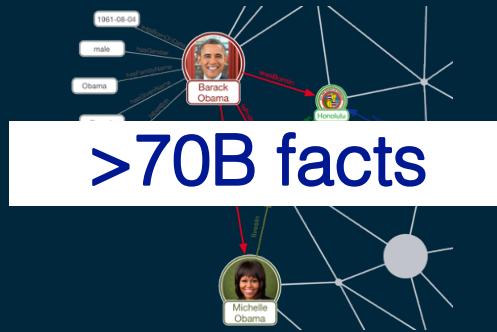
If you're interested in joining our group, send an email with your interests and CV to opportunities@umich.edu.



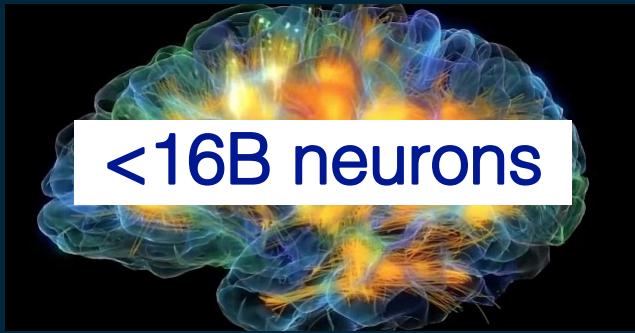
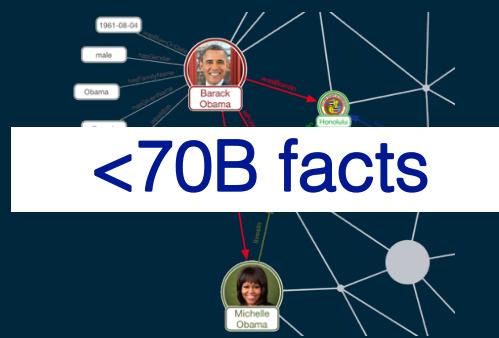
Graphs are everywhere!



LARGE-scale Graph Data



SUMMARIZATION of Big Datasets is Crucial!



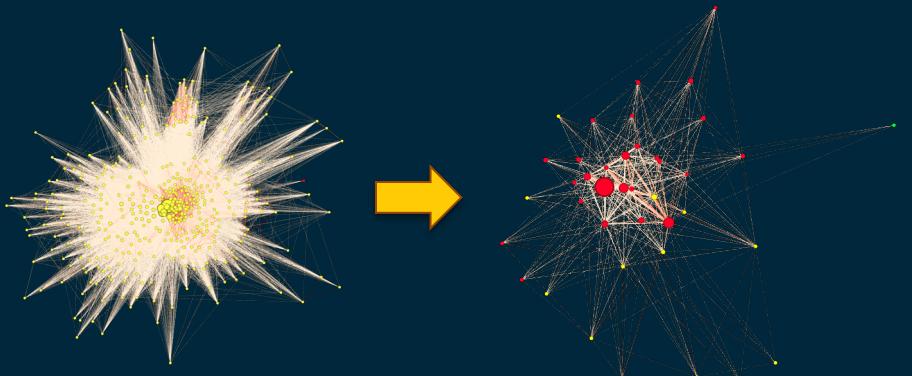
dblp



What is graph summarization?

Graph summarization seeks to find:

- a short representation of the input graph,
 - ✧ often in the form of an aggregated or sparsified graph, or a set of structures
- which reveals patterns in the original data and preserves specific structural or other properties, depending on the application domain.



Graph Summarization Methods and Applications: A Survey

YIKE LIU, TARA SAFAVI, ABHILASH DIGHE, and DANAI KOUTRA, University of Michigan, Ann Arbor

While advances in computing resources have made processing enormous amounts of data possible, human ability to identify patterns in such data has not scaled accordingly. Efficient computational methods for condensing and simplifying data are thus becoming vital for extracting actionable insights. In particular, while data summarization techniques have been studied extensively, only recently has summarizing interconnected data, or *graphs*, become popular. This survey is a structured, comprehensive overview of the state-of-the-art methods for summarizing graph data. We first broach the motivation behind and the challenges of graph summarization. We then categorize summarization approaches by the type of graphs taken as input and further organize each category by core methodology. Finally, we discuss applications of summarization on real-world graphs and conclude by describing some open problems in the field.

CCS Concepts: • Mathematics of computing → Graph algorithms; • Information systems → Data mining; Summarization; • Human-centered computing → Social network analysis; • Theory of computation → Unsupervised learning and clustering; • Computing methodologies → Network science;

Additional Key Words and Phrases: Graph mining, graph summarization

ACM Reference format:

Yike Liu, Tara Safavi, Abhilash Dighe, and Danai Koutra. 2018. Graph Summarization Methods and Applications: A Survey. *ACM Comput. Surv.* 51, 3, Article 62 (June 2018), 34 pages.
<https://doi.org/10.1145/3186727>

1 INTRODUCTION

As technology advances, the amount of data that we generate and our ability to collect and archive such data both increase continuously. Daily activities like social media interaction, web browsing, product and service purchases, itineraries, and wellness sensors generate large amounts of data, the analysis of which can immediately impact our lives. This abundance of generated data and its velocity call for data summarization, one of the main data mining tasks.

Since summarization facilitates the identification of structure and meaning in data, the data mining community has taken a strong interest in the task. Methods for a variety of data types

Y. Liu and T. Safavi contributed equally to this article.

This material was based on work supported in part by the National Science Foundation under grant IIS 1743088, Trove, and the University of Michigan. 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 National Science Foundation or other funding parties. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

Authors' addresses: Y. Liu, T. Safavi, A. Dighe, and D. Koutra, Bob and Betty Beyster Building, 2260 Hayward St, Ann Arbor, MI 48109; emails: yikeliu,tsafavi,adighe,dkoutra@umich.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2018 ACM 0360-0300/2018/06-ART62 \$15.00
<https://doi.org/10.1145/3186727>

ACM Computing Surveys, Vol. 51, No. 3, Article 62. Publication date: June 2018.

Why graph summarization?

- Reduction of data volume + storage
 - ✧ e.g., fewer I/O operations
- Speedup of algorithms + queries
- Interactive analysis
- Influence analysis and understanding
- Noise elimination -> reveals patterns
- Privacy preservation



Graph Summarization Methods and Applications: A Survey

YIKE LIU, TARA SAFAVI, ABHILASH DIGHE, and DANAI KOUTRA, University of Michigan, Ann Arbor

While advances in computing resources have made processing enormous amounts of data possible, human ability to identify patterns in such data has not scaled accordingly. Efficient computational methods for condensing and simplifying data are thus becoming vital for extracting actionable insights. In particular, while data summarization techniques have been studied extensively, only recently has summarizing interconnected data, or *graphs*, become popular. This survey is a structured, comprehensive overview of the state-of-the-art methods for summarizing graph data. We first broach the motivation behind and the challenges of graph summarization. We then categorize summarization approaches by the type of graphs taken as input and further organize each category by core methodology. Finally, we discuss applications of summarization on real-world graphs and conclude by describing some open problems in the field.

CCS Concepts: • Mathematics of computing → Graph algorithms; • Information systems → Data mining; Summarization; • Human-centered computing → Social network analysis; • Theory of computation → Unsupervised learning and clustering; • Computing methodologies → Network science;

Additional Key Words and Phrases: Graph mining, graph summarization

ACM Reference format:

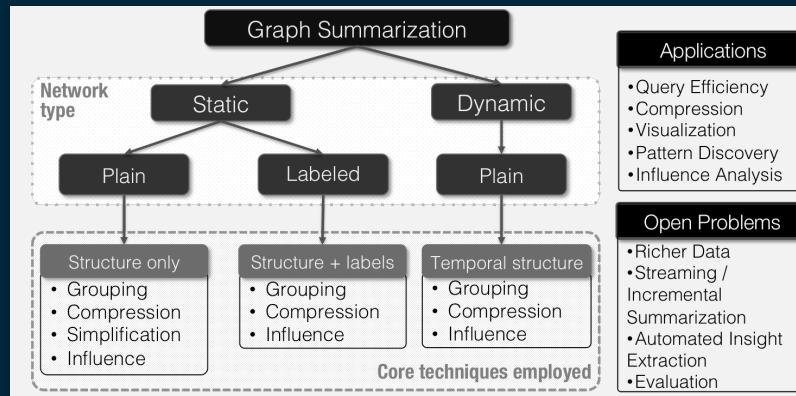
Yike Liu, Tara Safavi, Abhilash Dighe, and Danai Koutra. 2018. Graph Summarization Methods and Applications: A Survey. *ACM Comput. Surv.* 51, 3, Article 62 (June 2018), 34 pages.
<https://doi.org/10.1145/3186727>

1 INTRODUCTION

As technology advances, the amount of data that we generate and our ability to collect and archive such data both increase continuously. Daily activities like social media interaction, web browsing, product and service purchases, itineraries, and wellness sensors generate large amounts of data, the analysis of which can immediately impact our lives. This abundance of generated data and its velocity call for data summarization, one of the main data mining tasks.

Since summarization facilitates the identification of structure and meaning in data, the data mining community has taken a strong interest in the task. Methods for a variety of data types

Y. Liu and T. Safavi contributed equally to this article.
This material was based on work supported in part by the National Science Foundation under grant IIS 1743088, Trove, and the University of Michigan. 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 National Science Foundation or other funding parties.



Why graph summarization?

- Reduction of data volume + storage
 - ✧ e.g., fewer I/O operations
- Speedup of algorithms + queries
- Interactive analysis
- Influence analysis and understanding
- Noise elimination -> reveals patterns
- Privacy preservation + anomalies



Graph Summarization Methods and Applications: A Survey

YIKE LIU, TARA SAFAVI, ABHILASH DIGHE, and DANAI KOUTRA, University of Michigan, Ann Arbor

While advances in computing resources have made processing enormous amounts of data possible, human ability to identify patterns in such data has not scaled accordingly. Efficient computational methods for condensing and simplifying data are thus becoming vital for extracting actionable insights. In particular, while data summarization techniques have been studied extensively, only recently has summarizing interconnected data, or *graphs*, become popular. This survey is a structured, comprehensive overview of the state-of-the-art methods for summarizing graph data. We first broach the motivation behind and the challenges of graph summarization. We then categorize summarization approaches by the type of graphs taken as input and further organize each category by core methodology. Finally, we discuss applications of summarization on real-world graphs and conclude by describing some open problems in the field.

CCS Concepts: • Mathematics of computing → Graph algorithms; • Information systems → Data mining; Summarization; • Human-centered computing → Social network analysis; • Theory of computation → Unsupervised learning and clustering; • Computing methodologies → Network science;

Additional Key Words and Phrases: Graph mining, graph summarization

ACM Reference format:

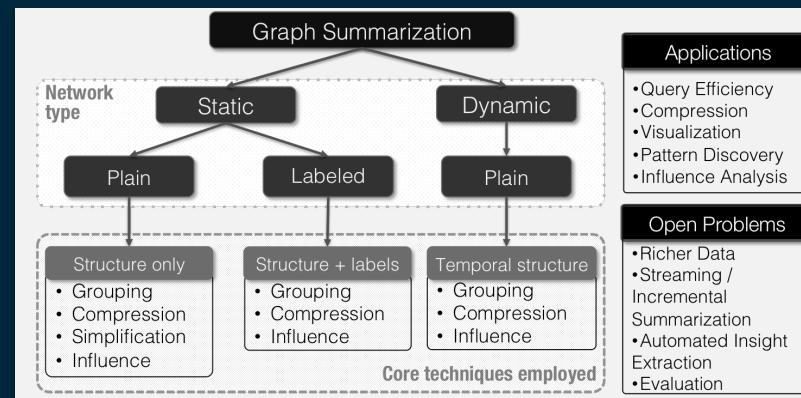
Yike Liu, Tara Safavi, Abhilash Dighe, and Danai Koutra. 2018. Graph Summarization Methods and Applications: A Survey. *ACM Comput. Surv.* 51, 3, Article 62 (June 2018), 34 pages.
<https://doi.org/10.1145/3186727>

1 INTRODUCTION

As technology advances, the amount of data that we generate and our ability to collect and archive such data both increase continuously. Daily activities like social media interaction, web browsing, product and service purchases, itineraries, and wellness sensors generate large amounts of data, the analysis of which can immediately impact our lives. This abundance of generated data and its velocity call for data summarization, one of the main data mining tasks.

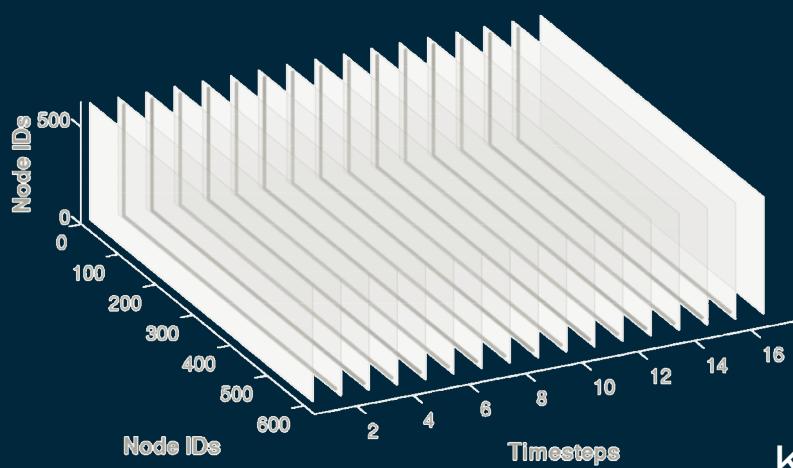
Since summarization facilitates the identification of structure and meaning in data, the data mining community has taken a strong interest in the task. Methods for a variety of data types

Y. Liu and T. Safavi contributed equally to this article.
This material was based on work supported in part by the National Science Foundation under grant IIS 1743088, Trove, and the University of Michigan. 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 National Science Foundation or other funding parties.



Summarizing Large Networks: Overview

Survey:
[CSUR'18]

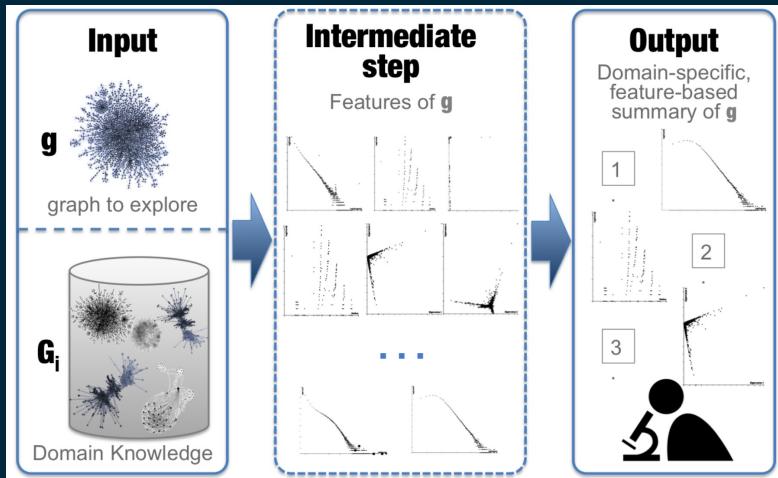


Structural Summaries

[SDM'14, KDD'15,
Dat Bull Eng'17,
SNAM'18,
SDM'19,
KDD'19a, KDD'20...]

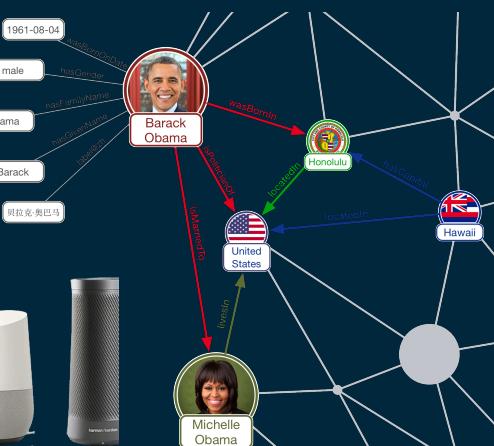
Domain-specific Summaries

[ICDM'17,
KDD'19]



Query-on-the-edge + Rule-based Summaries

[ICDM'19,
WebConf'20]

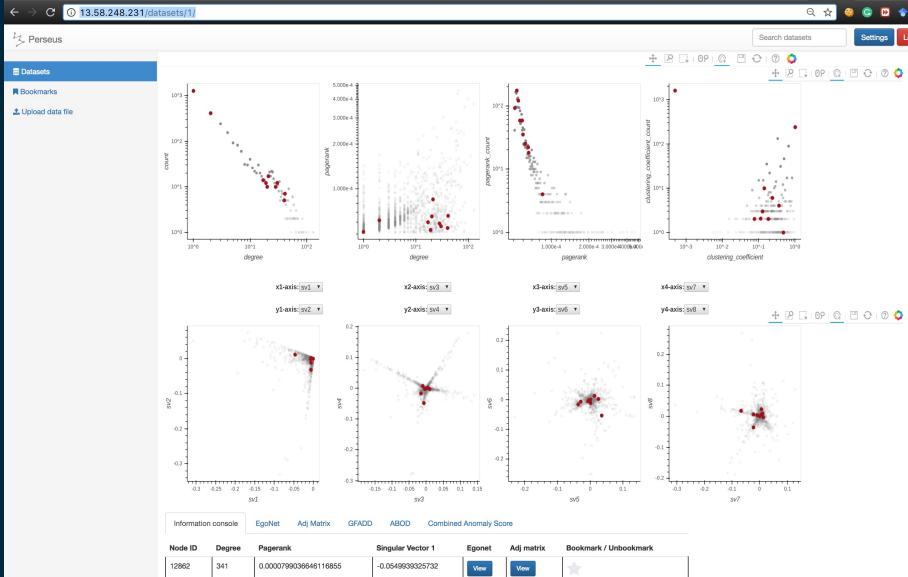


Interactive Summaries

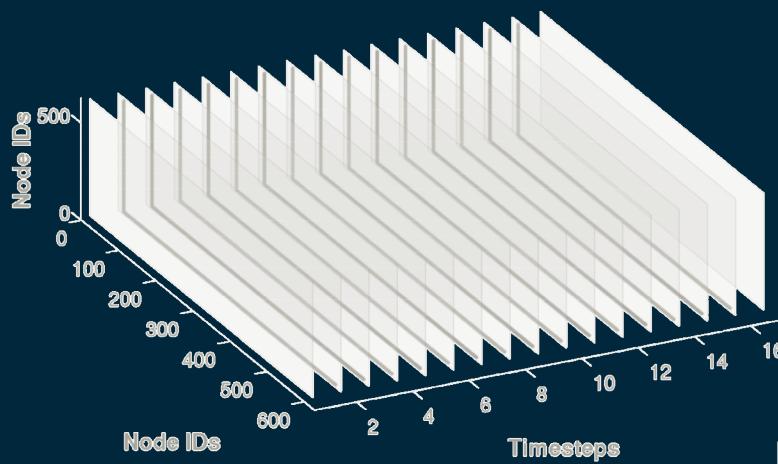
Latent Summaries [KDD'19a, b]

$$M = U V^T$$

[VLDB'15, Informatics'17]



This talk: Summarization Meets Outlier Detection



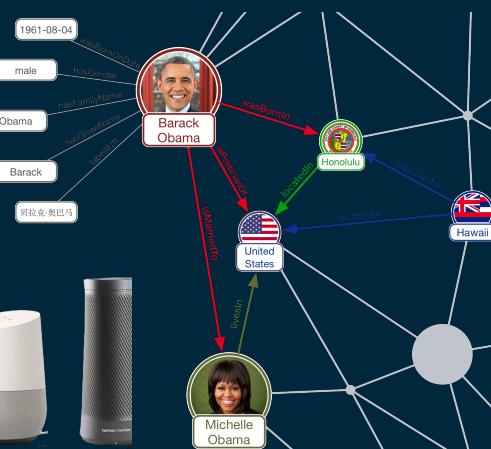
Structural Summaries

[SDM'14, KDD'15,
Dat Bull Eng'17,
SNAM'18,
SDM'19,
KDD'19a, KDD'20...]

Query-on-the-edge + Rule-based

Summaries

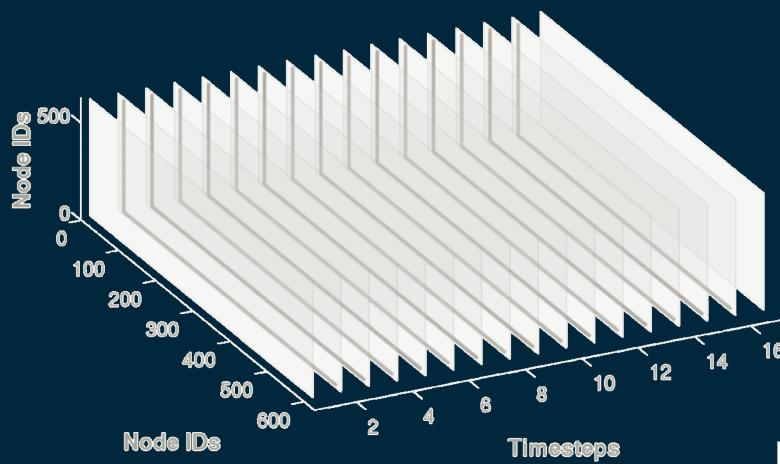
[ICDM'19,
WebConf'20]



② Graph Streams:
Persistent and bursty activity detection
[KDD'20]

① Knowledge Graphs:
Unified error detection and completion
[WebConf'20]

This talk: Summarization Meets Outlier Detection



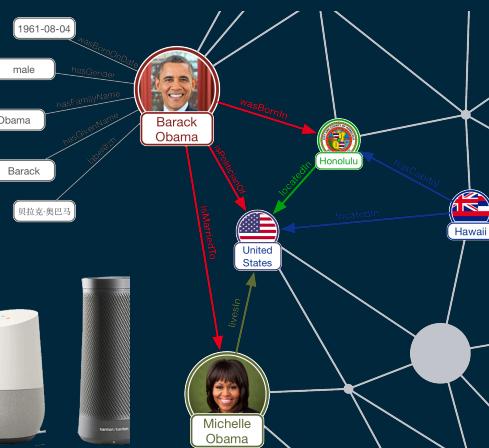
Structural Summaries

[SDM'14, KDD'15,
Dat Bull Eng'17,
SNAM'18,
SDM'19,
KDD'19a, KDD'20...]

Query-on-the-edge + Rule-based

Summaries

[ICDM'19,
WebConf'20]

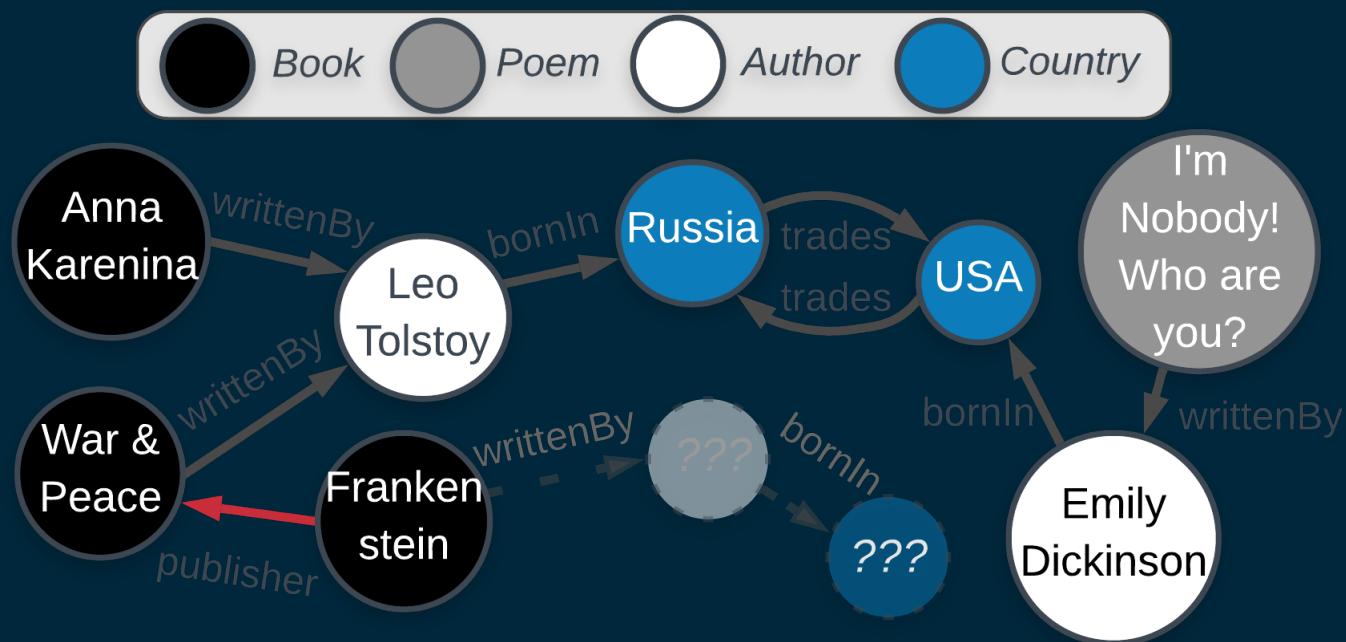


② Graph Streams:
Persistent and bursty activity detection
[KDD'20]

① Knowledge Graphs:
Unified error detection and completion
[WebConf'20]

Knowledge graphs (KGs)

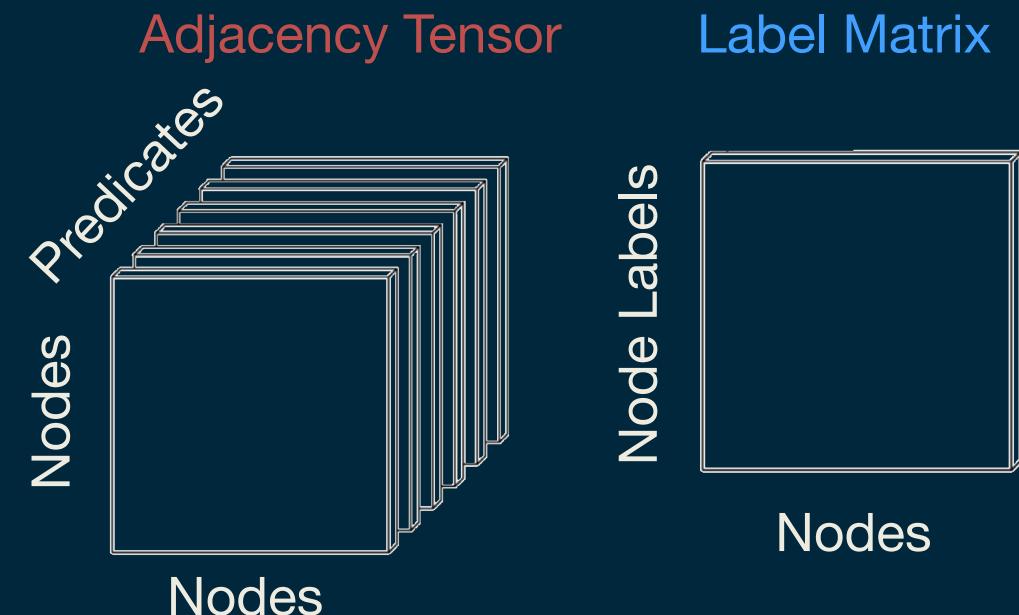
store general information about the world in the structure of a graph



edge = triple (subject node or head,
predicate or relation,
object node or tail)

Knowledge graphs (KGs)

can be represented as labeled, directed multi-relational graphs



practice

Article development led by ACM Queue
queue.acm.org

DOI:10.1145/3331166

**Five diverse technology companies
show how it's done.**

BY NATASHA NOY, YUQING GAO, ANSHU JAIN,
ANANT NARAYANAN, ALAN PATTERSON, AND JAMIE TAYLOR

Industry-Scale Knowledge Graphs: Lessons and Challenges

KNOWLEDGE GRAPHS ARE critical to many enterprises today: They provide the structured data and factual knowledge that drive many products and make them more intelligent and “magical.”

In general, a knowledge graph describes objects of interest and connections between them. For example, a knowledge graph may have nodes for a movie, the actors in this movie, the director, and so on. Each node may have properties such as an actor’s name and age. There may be nodes for multiple movies involving a particular actor. The user can then traverse the knowledge graph to collect information on all the movies in which the actor appeared or, if applicable, directed.

Many practical implementations impose constraints on the links in knowledge graphs by defining a *schema* or *ontology*. For example, a link from a movie to its director must connect an object of type Movie to an object of type Person. In some cases the links themselves might have their own properties: a link connecting an actor and a movie might have the name of the specific role the actor played. Similarly, a link connecting a politician with a specific role in government might have the time period during which the politician held that role.

Knowledge graphs and similar structures usually provide a shared substrate of knowledge within an organization, allowing different products and applications to use similar vocabulary and to reuse definitions and descriptions that others create. Furthermore, they usually provide a compact formal representation that developers can use to infer new facts and build up the knowledge—for example, using the graph connecting movies and actors to find out which actors frequently appear in movies together.

This article looks at the knowledge graphs of five diverse tech companies, comparing the similarities and differences in their respective experiences of building and using the graphs, and discussing the challenges that all knowledge-driven enterprises face today. The collection of knowledge graphs discussed here covers the breadth of applications, from search, to product descriptions, to social networks:

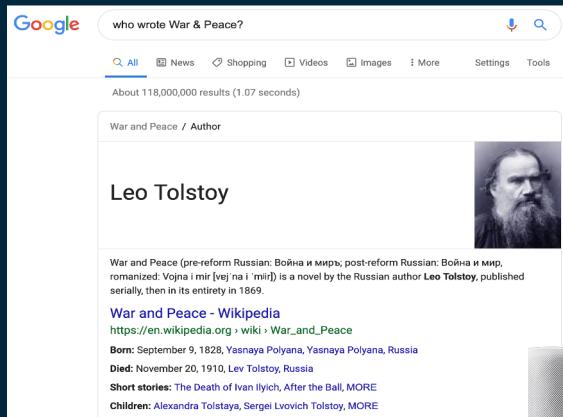
- Both Microsoft’s Bing knowledge graph and the Google Knowledge Graph support search and answering questions in search and during conversations. Starting with the descriptions and connections of people, places, things, and organizations, these graphs include general knowledge about the world.

- Facebook has the world’s largest social graph, which also includes information about music, movies, celebrities, and places that Facebook users care about.



Applications of KGs

Question Answering & Chatbots



Automatic Fact Checking



Reading Comprehension



Semantic search
Biomedical applications

...

Financial applications
Recommendation systems

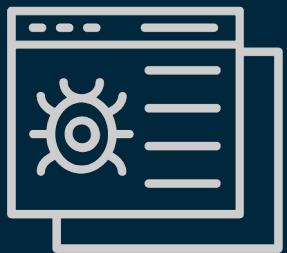
...

KGs are constructed via

Crowd
Sourcing

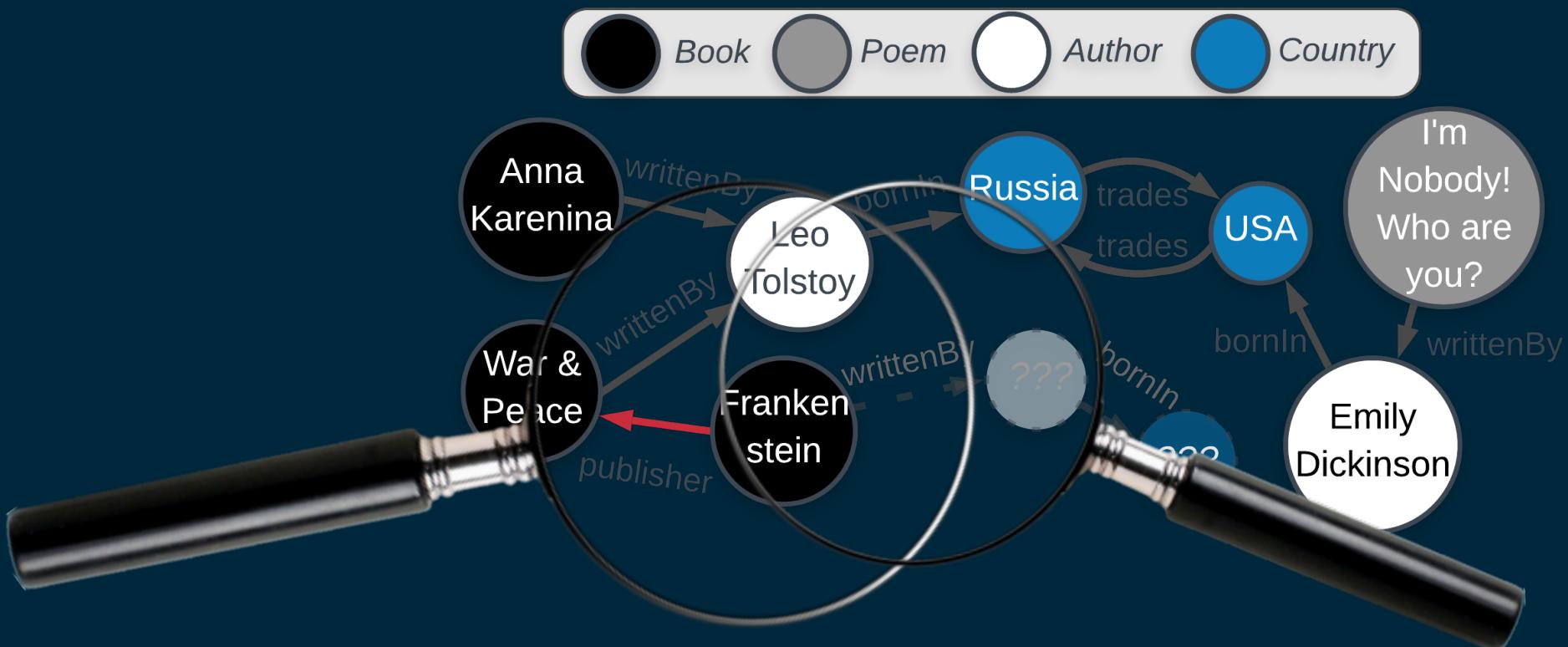


Web
Crawling

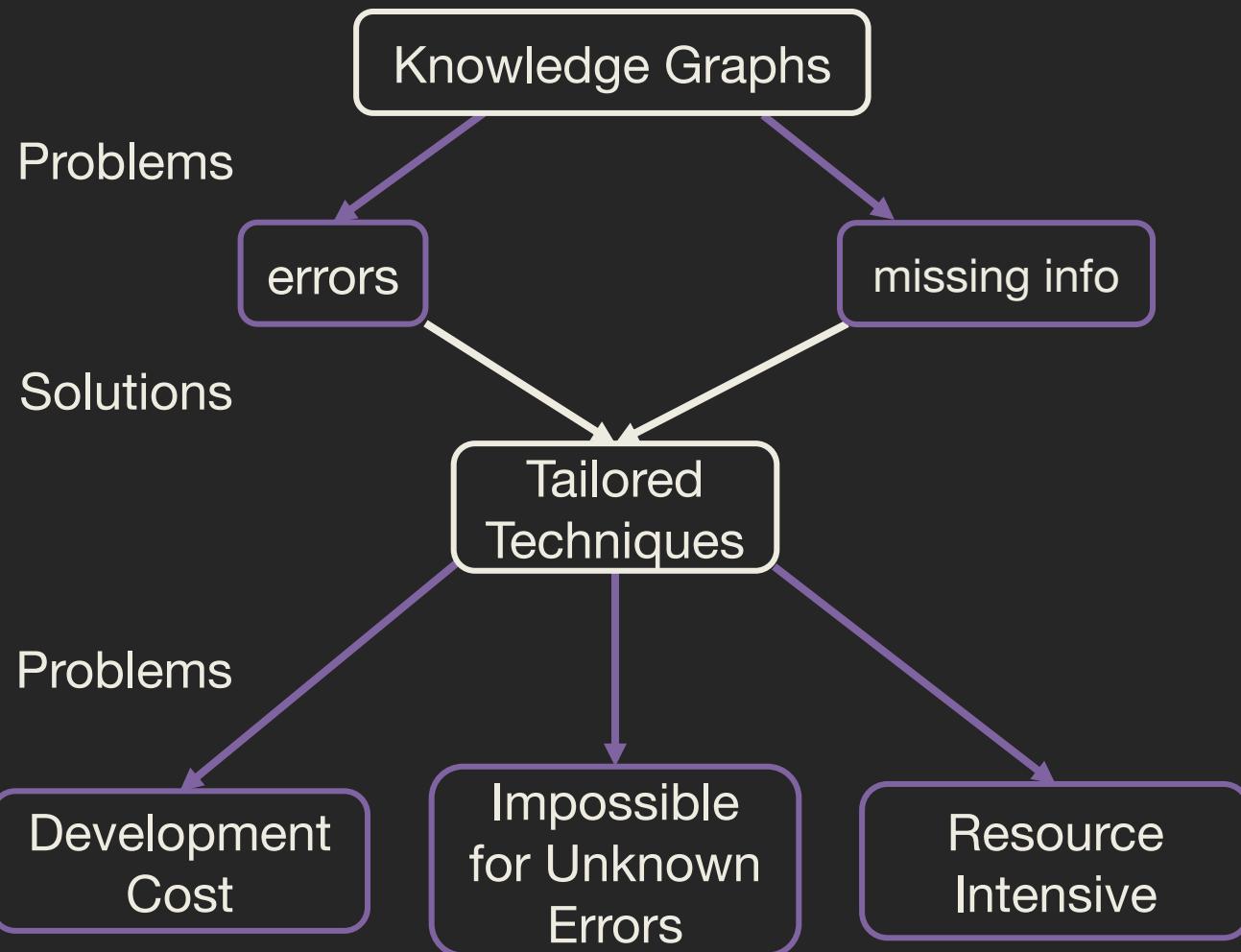


...which leads to

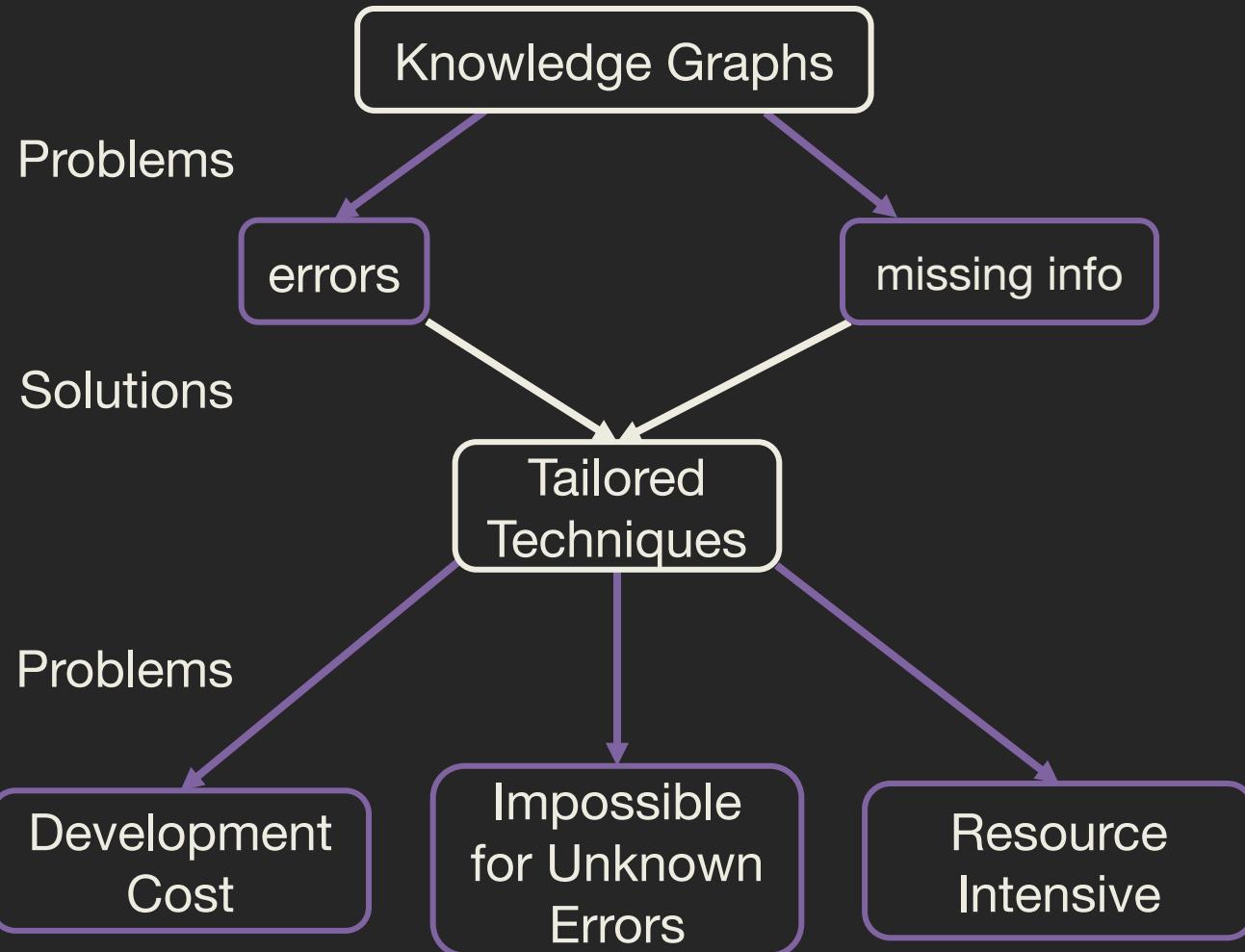
errors and miss some information



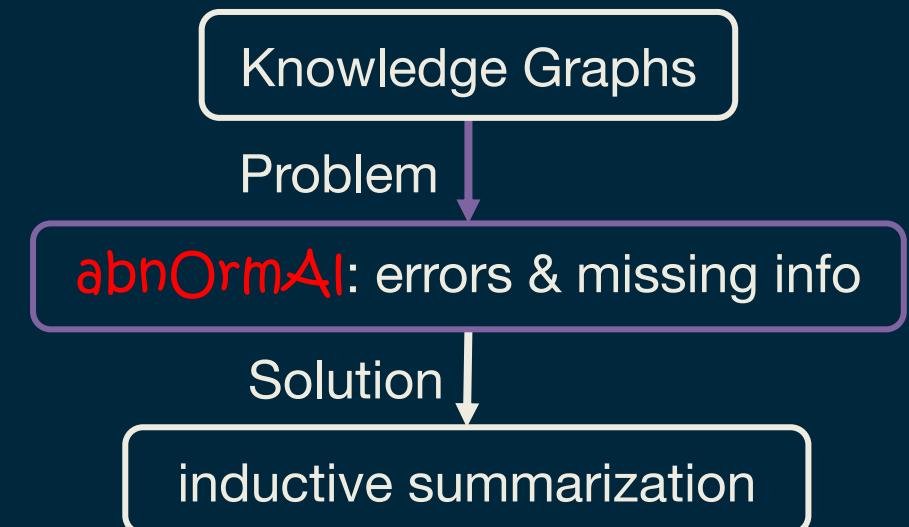
Current Approach



Current Approach



Proposed Approach



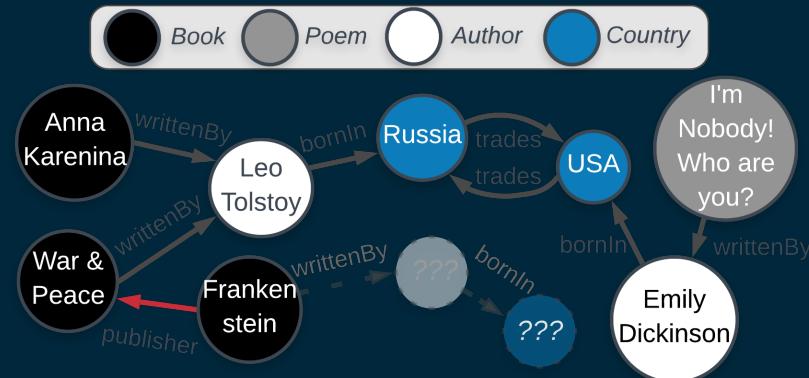
KGIST: Knowledge Graph Inductive SummarizaTion

Find a concise summary M of knowledge graph G ,
consisting of inductive, soft rules s.t.

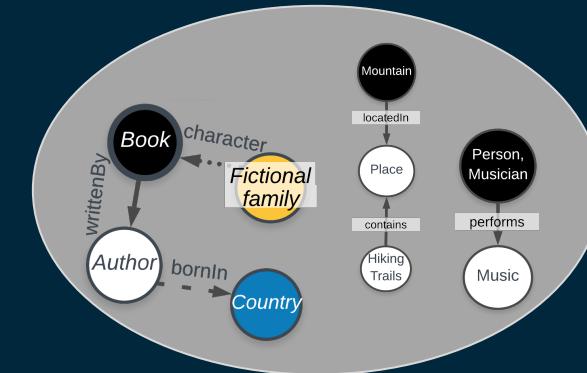
$$\min L(G, M) = L(M) + L(G|M)$$

bits to describe M

bits to describe G with M



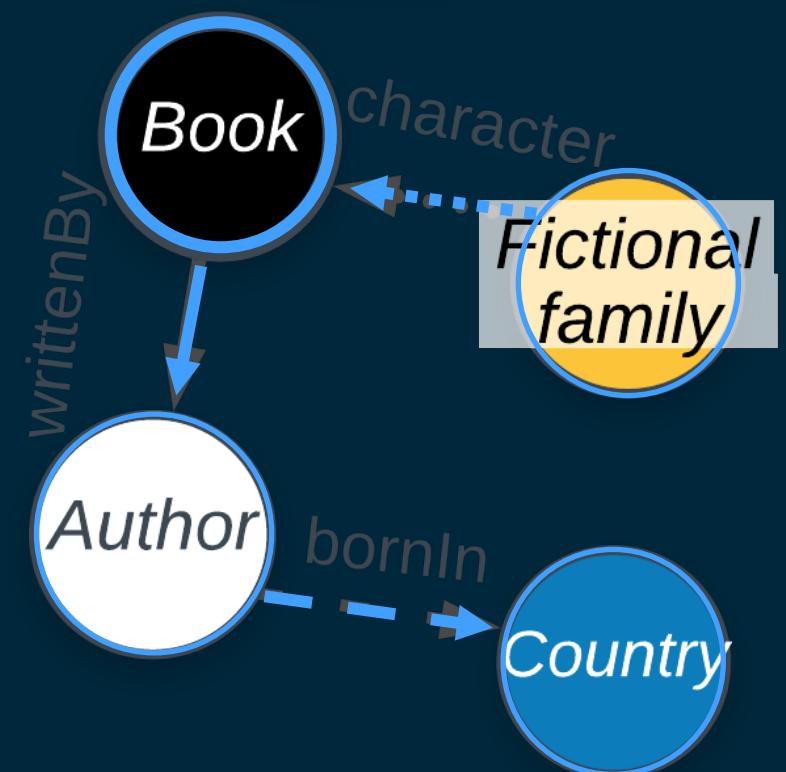
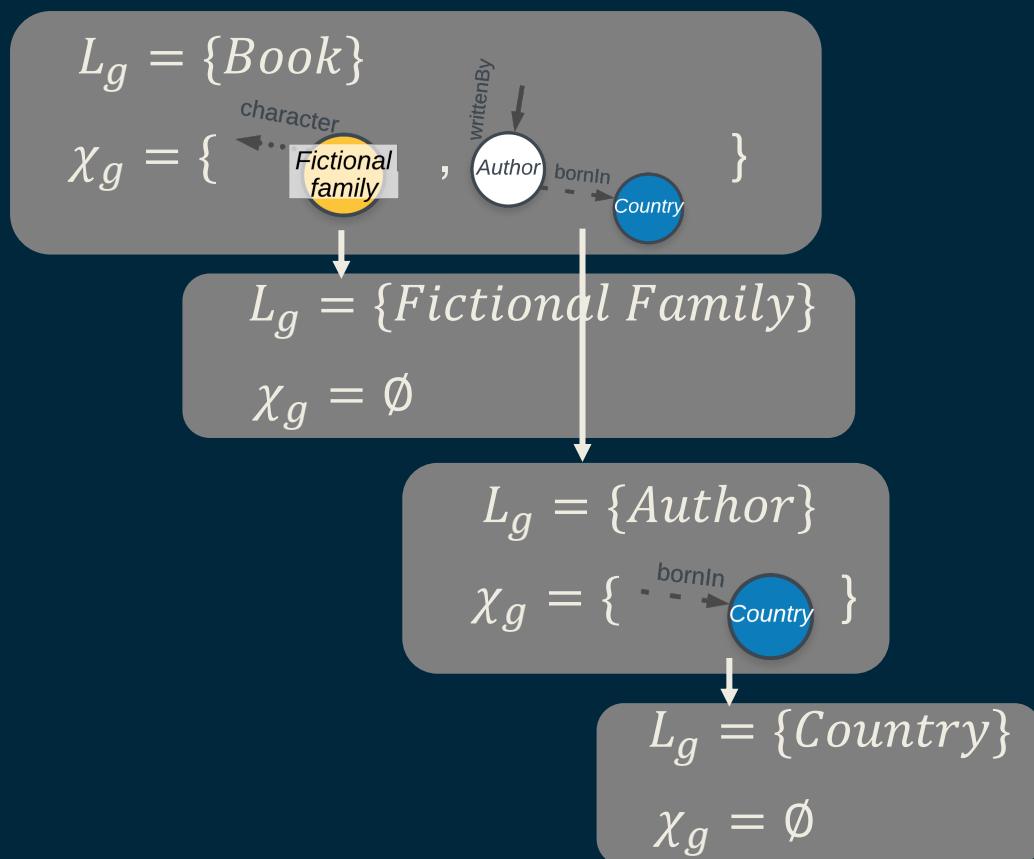
$M =$



Rule $g = (L_g, \chi_g)$ = (root label, children rules)

We formulate rules **recursively** as rooted, directed, and labeled graphs

- A rule asserts things about nodes with the root labels, L_g

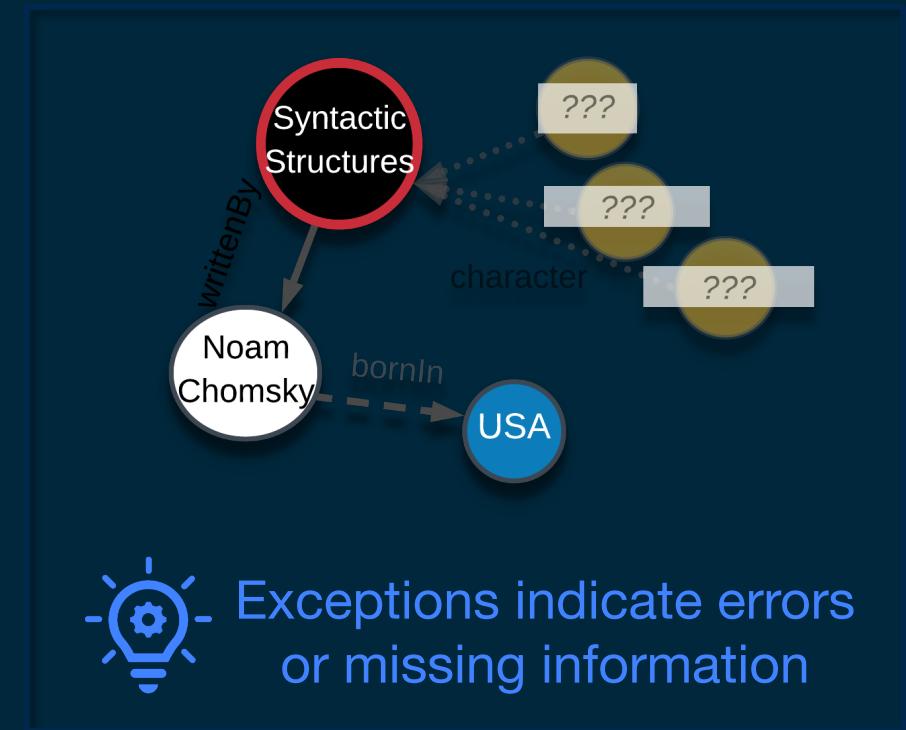


Correct assertions $\mathcal{A}_c^{(g)}$ & Exceptions $\mathcal{A}_{\xi}^{(g)}$

Guided traversals
that a rule implies



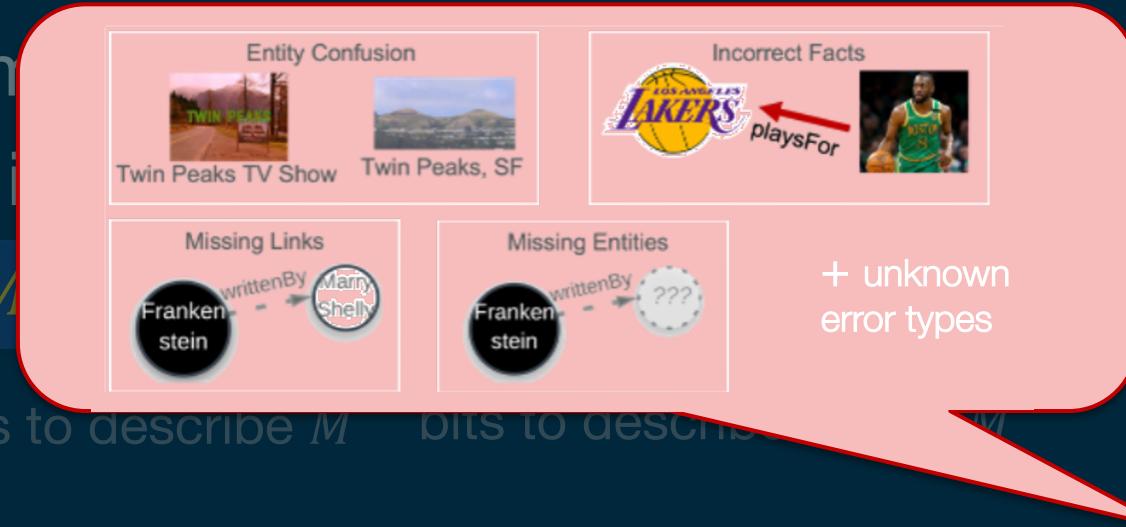
Rule g



KGIST: Knowledge Graph Inductive Summarization

Find a concise summary
consisting of

$$\min L(G, M)$$

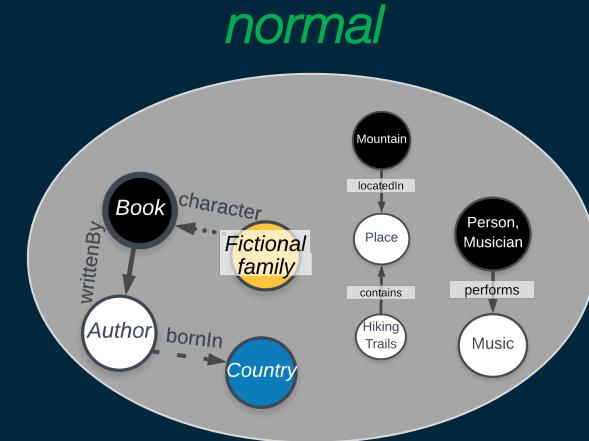


bits to describe M

bits to describe M

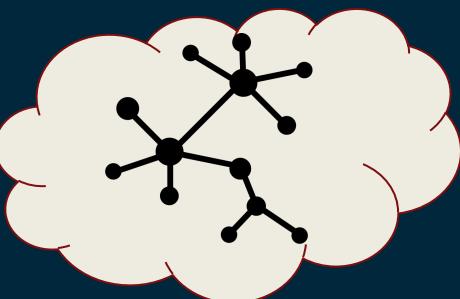


$$M =$$



Deriving $L(G, M) = L(M) + L(G|M)$

Alice (sender)



Hey Alice, could
you tell me
about your KG?

Bob (receiver)



MDL Model: Overview

Alice



Sure! I'll send:

- 1) Model-independent information
- 2) A model M
- 3) Any error the model makes

Bob

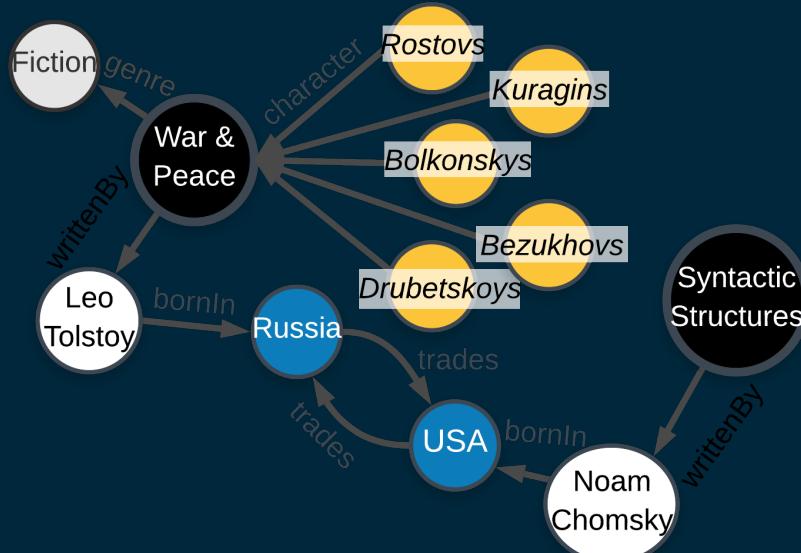


Ok, send the model the minimizes

$$L(G, M) = L(M) + L(G | M)$$



$$\text{MDL Model: } L(G, M) = L(M) + L(G|M)$$



Alice



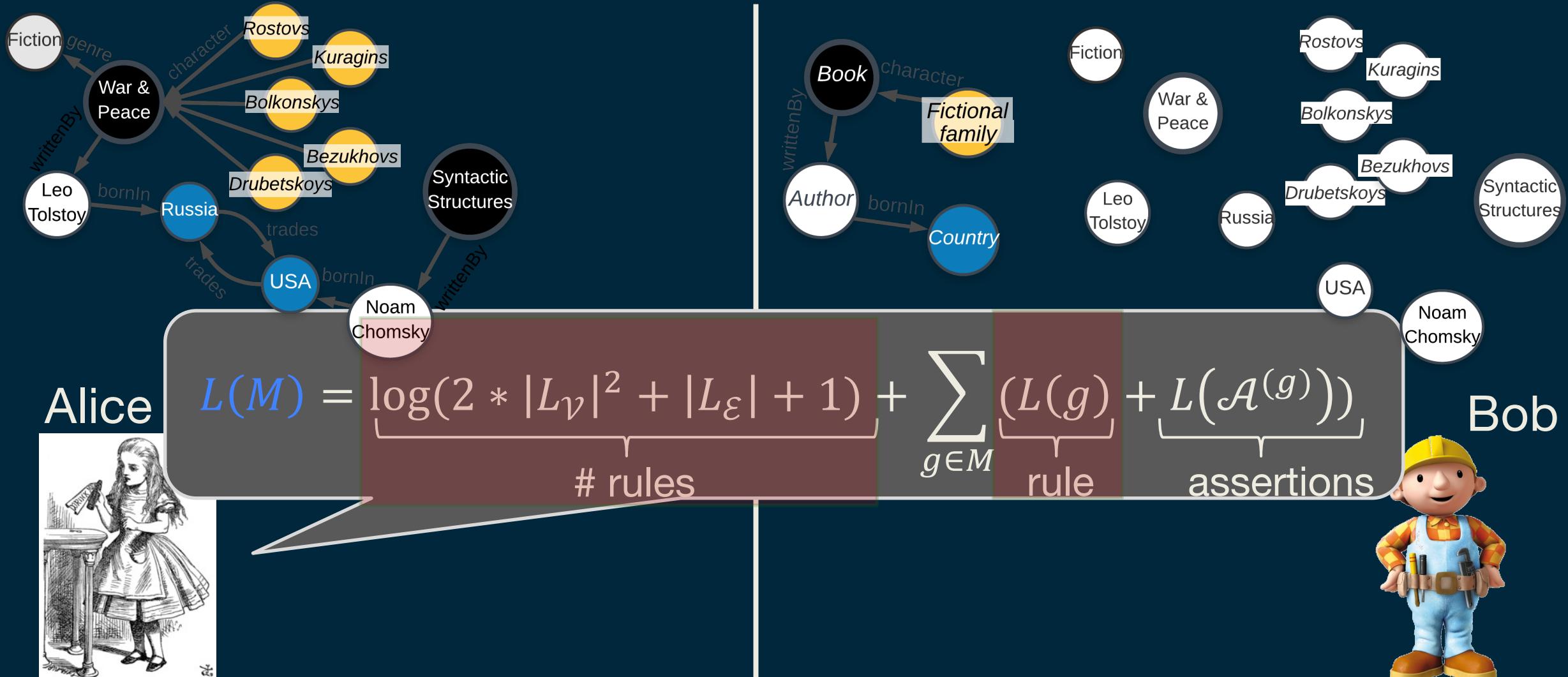
Model independent info:
nodes, # edges, node ids ...



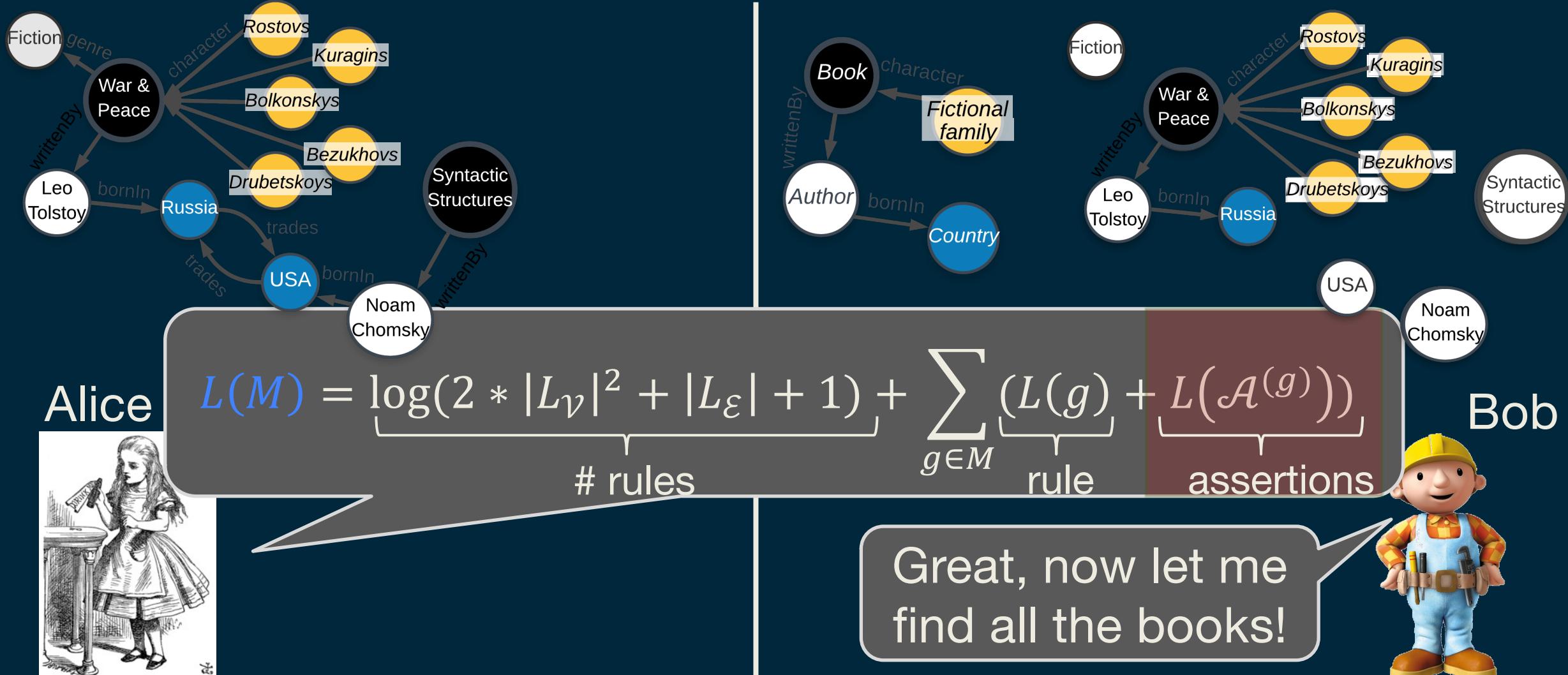
Bob



MDL Model: $L(G, M) = L(M) + L(G|M)$



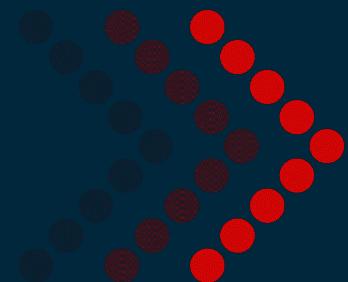
MDL Model: $L(G, M) = L(M) + L(G|M)$



MDL Model: $L(G, M) = L(M) + L(G|M)$

- Alice continues with the assertions, traversals etc...
- Done with the definition of $L(M)$

Alice



$\forall \exists \pi \cdot \infty$



$\forall \exists \pi \cdot \infty$



Bob



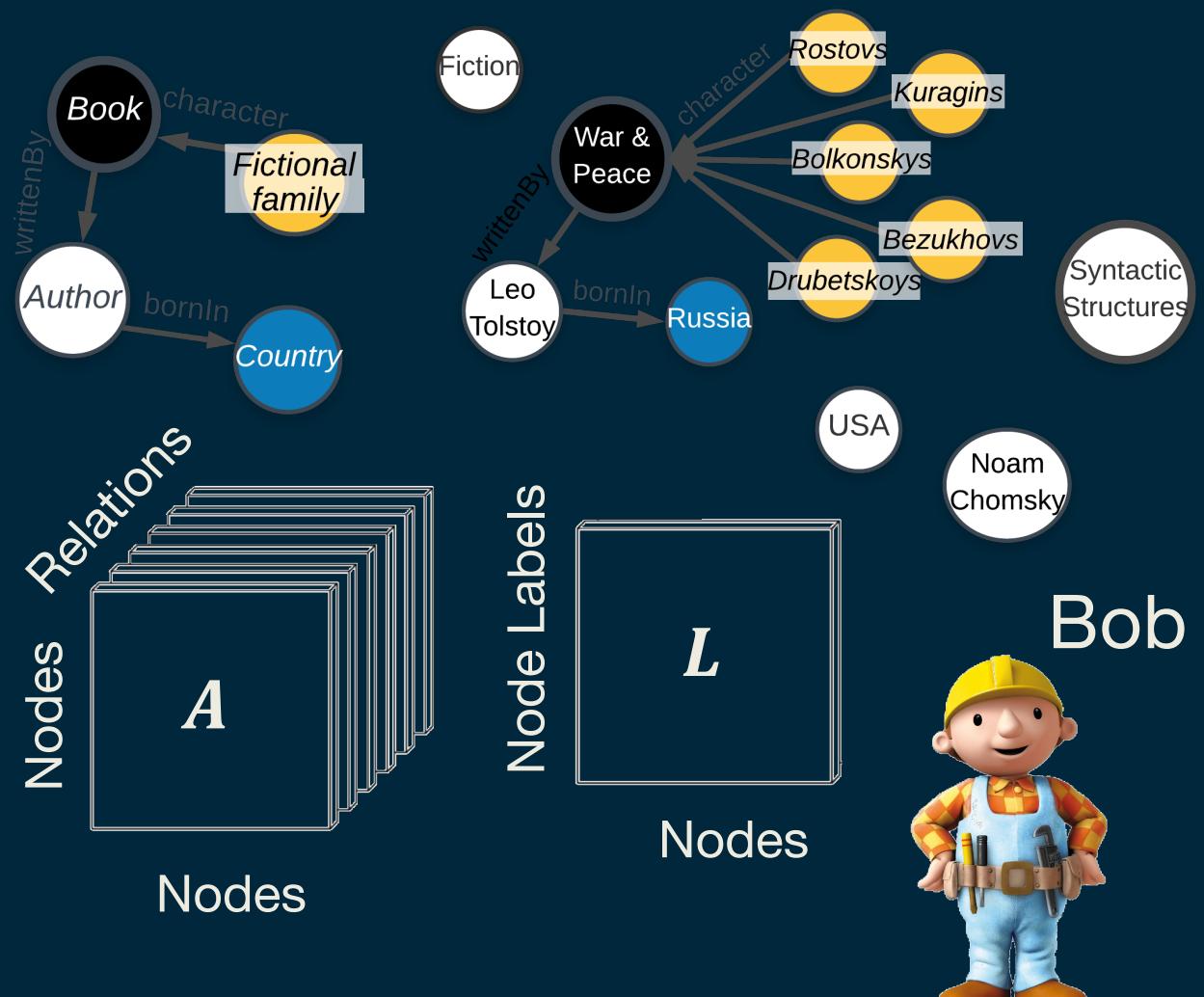
$$\text{MDL Model: } L(G, M) = L(M) + L(G|M)$$

Alice



$$L(G|M) = L(L^-) + L(A^-)$$

I'll send the 1s in
 L and A that the
rules didn't reveal

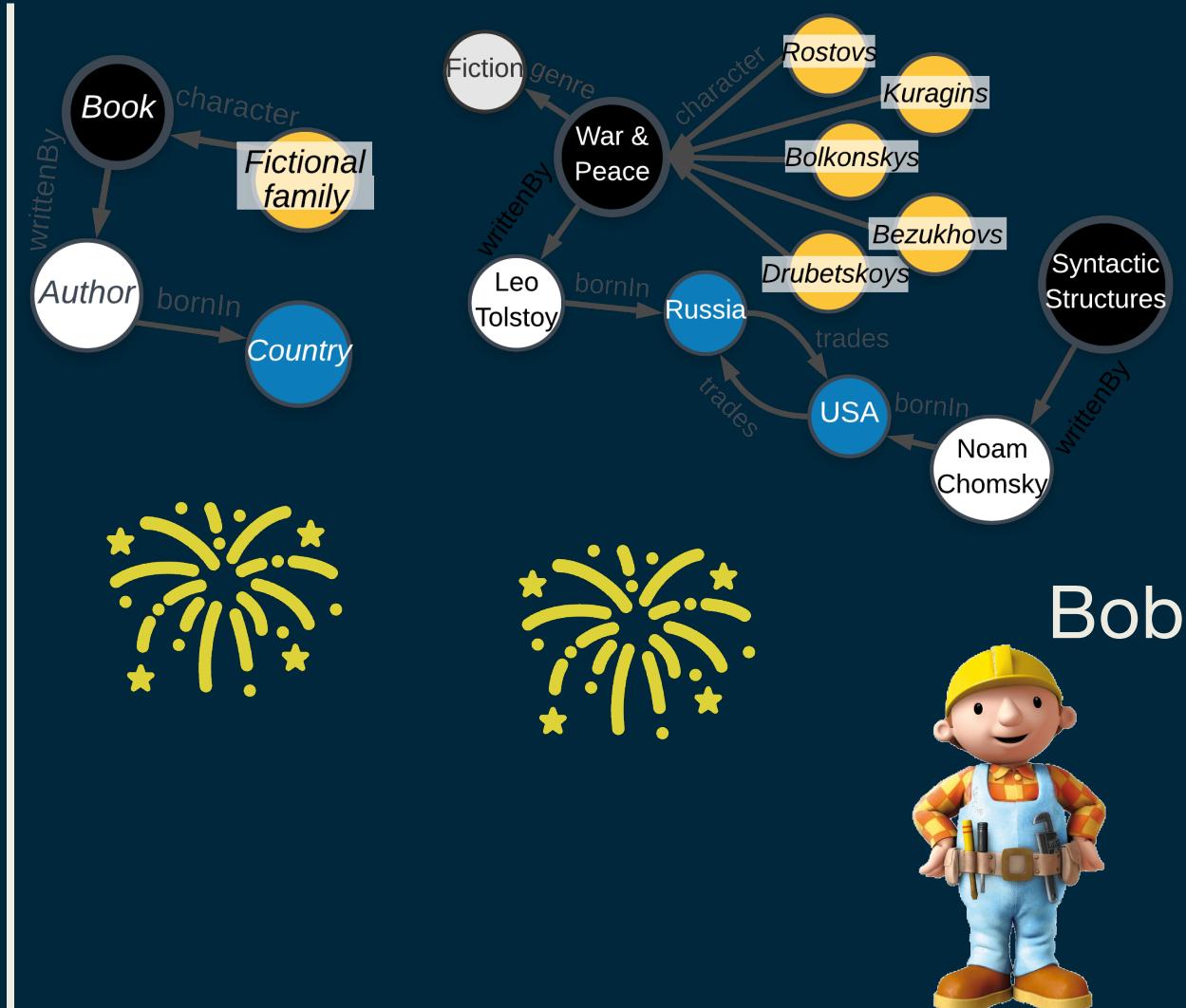


$$\text{MDL Model: } L(G, M) = L(M) + L(G|M)$$



Alice

There you go!



KGIST: Knowledge Graph Inductive Summarization

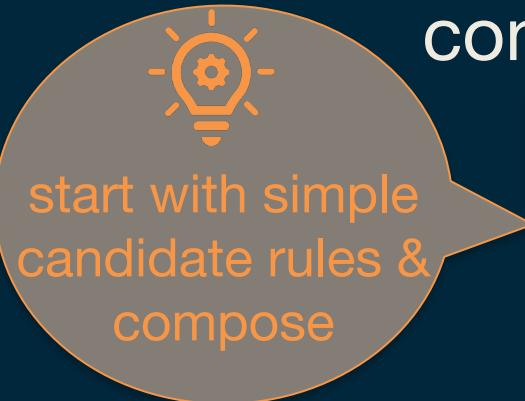
Find a concise summary M of knowledge graph G ,
consisting of inductive, soft rules s.t.

$$\min L(G, M) = L(M) + L(G|M)$$

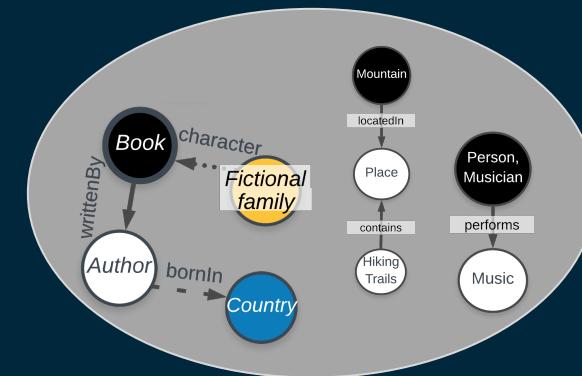
?

bits to describe M

bits to describe G with M

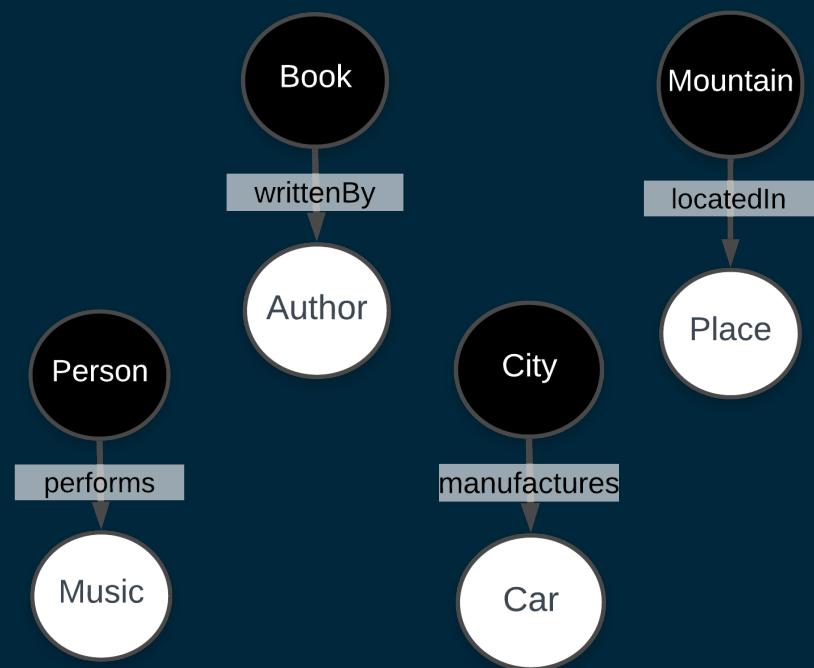


$M =$



KGIST Method: Overview

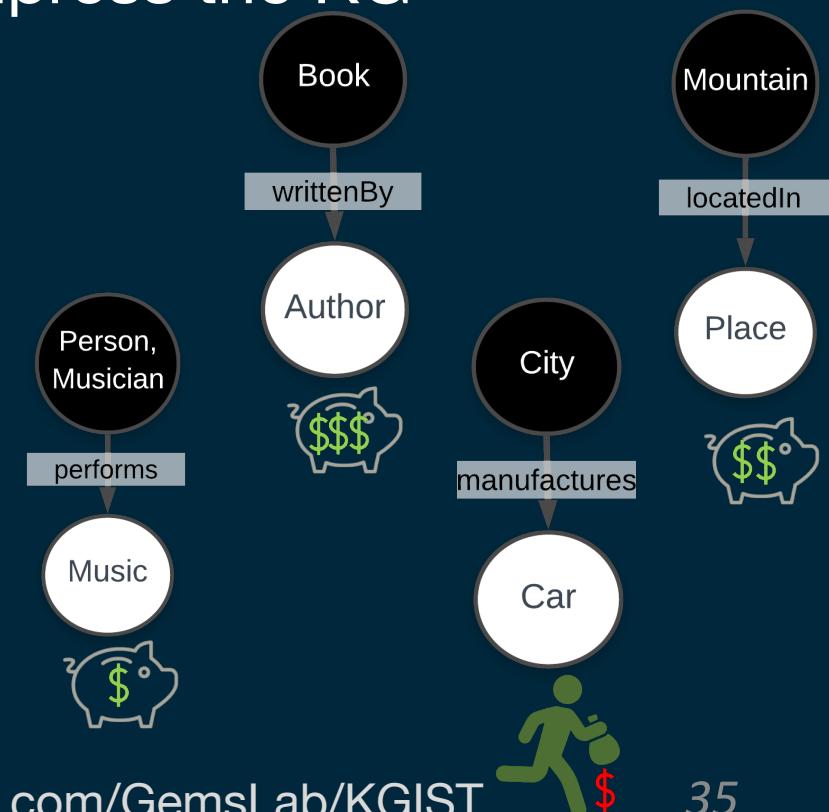
1. Generate candidate rules



KGIST Method: Overview

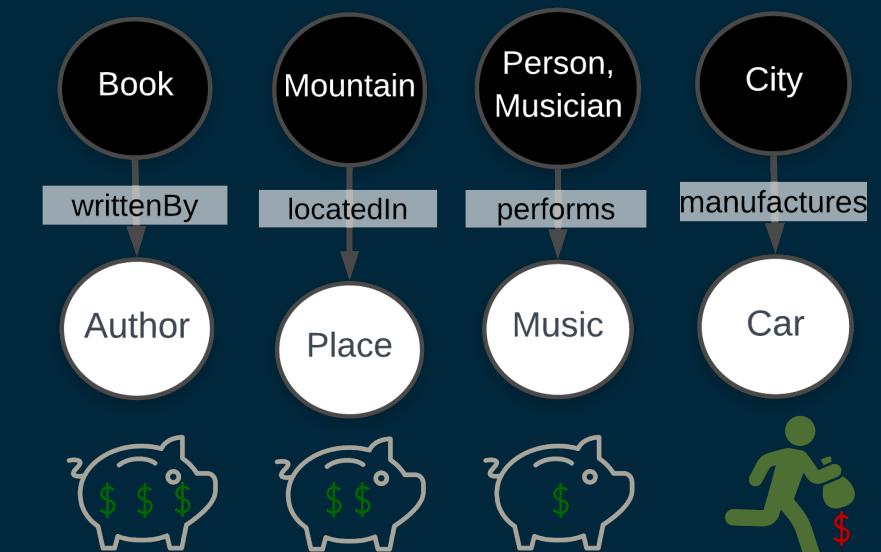
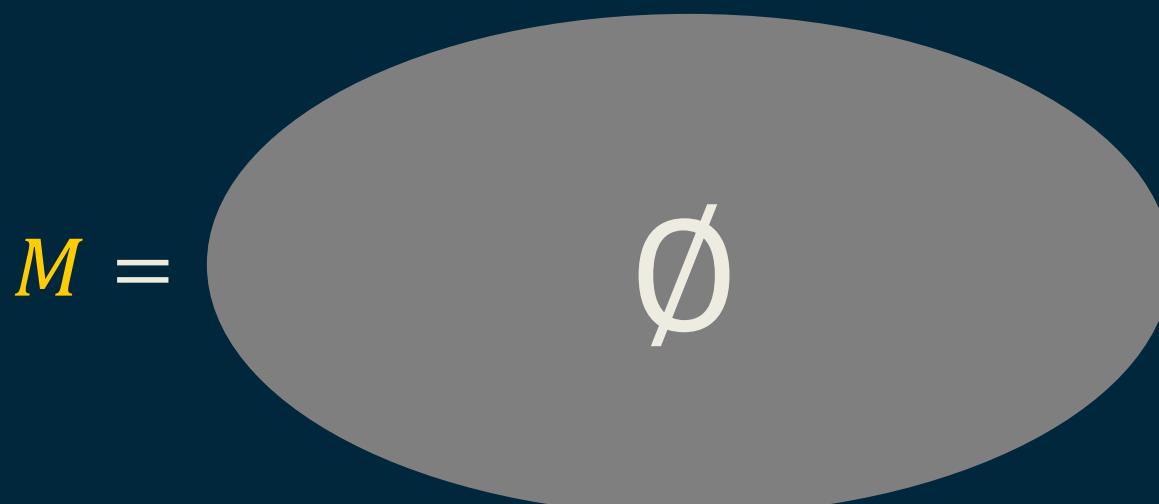
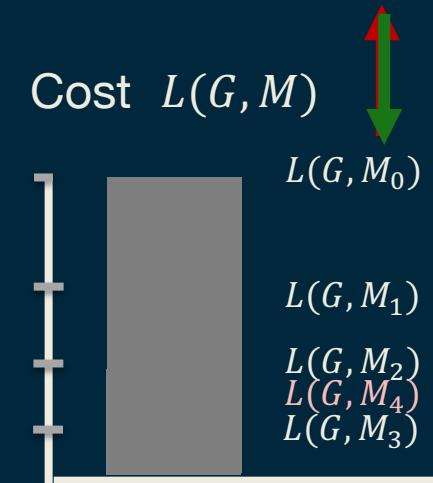
1. Generate candidate rules
2. Rank candidate rules
 - ✧ Based on how much they help explain/compress the KG

$$\underbrace{\Delta L(G|M_0 \cup \{g\})}_{\text{how much } g \text{ explains}} = \underbrace{L(G|M_0)}_{\# \text{ bits w/o } g} - \underbrace{L(G|M_0 \cup \{g\})}_{\# \text{ bits with } g}$$



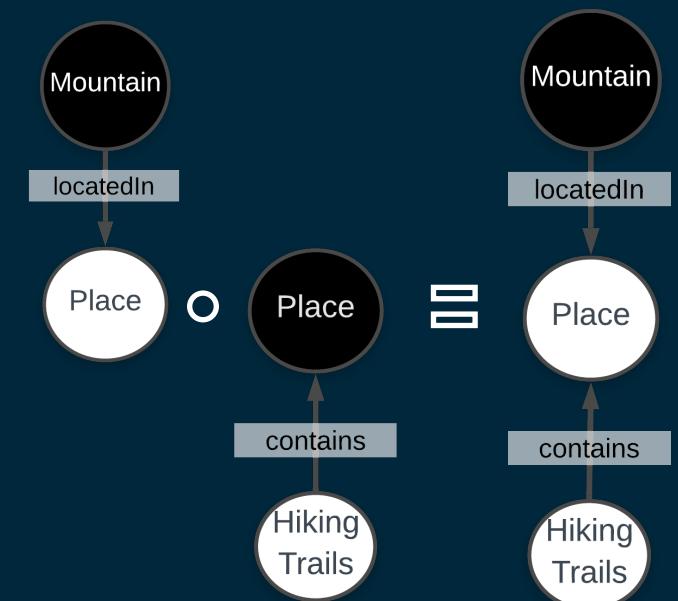
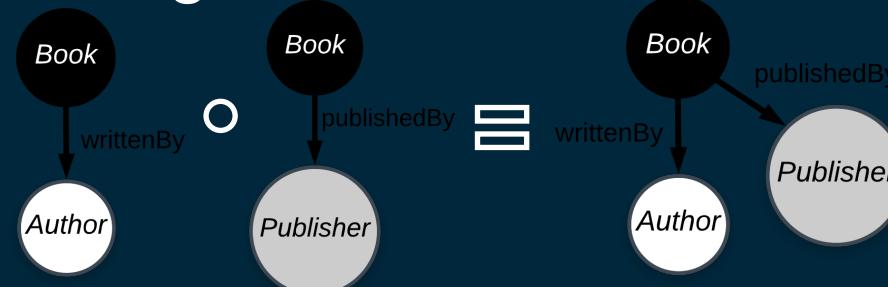
KGIST Method: Overview

1. Generate candidate rules
2. Rank candidate rules
 - ◊ Based on how much they help explain/compress the KG
3. Select rules
 - ◊ Based on minimizing $L(G, M)$

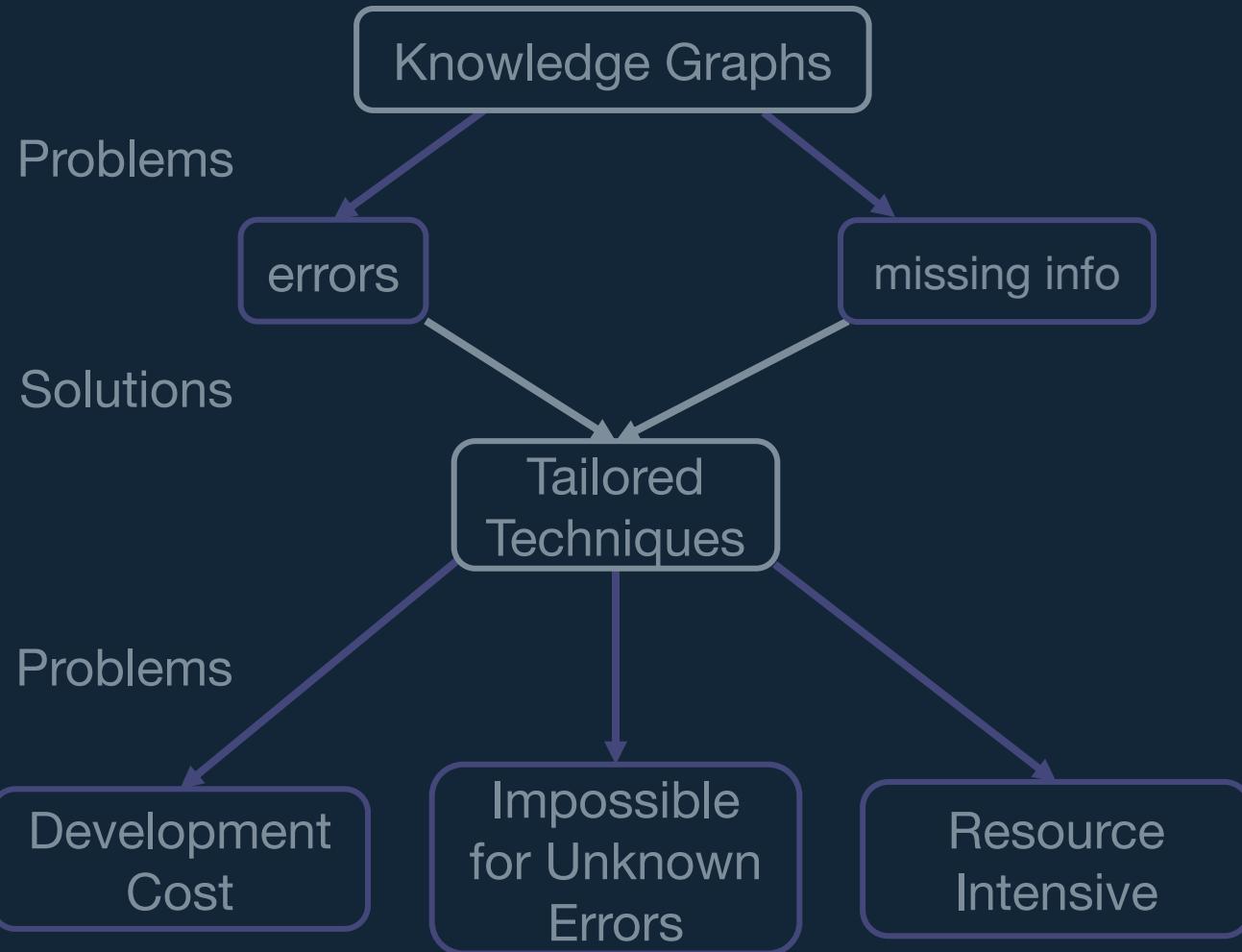


KGIST Method: Overview

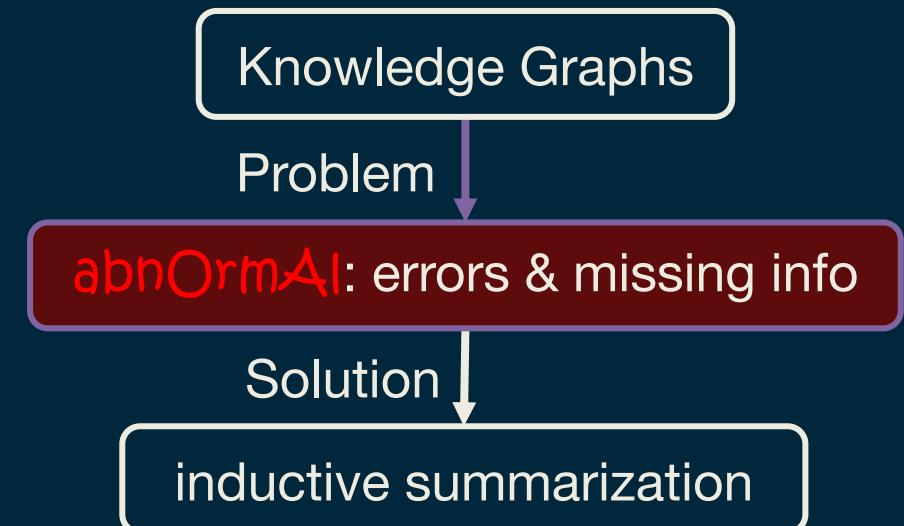
1. Generate candidate rules
2. Rank candidate rules
 - ✧ Based on how much they help explain/compress the KG
3. Select rules
 - ✧ Based on minimizing $L(G, M)$
4. Refine rules
 - ✧ Merging and nesting



Current Approach



Proposed Approach: KGIST



KGIST Anomaly Scores

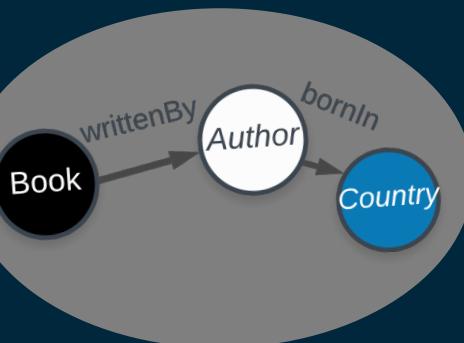
- Anomalous entities: violate many rules
 - MDL intuition: many bits to describe a node as an exception

$$\eta(v) = \sum_{\substack{g \in r(v): \\ v \in \mathcal{A}_\xi^{(g)}}} \frac{1}{|\mathcal{A}_\xi^{(g)}|} \log \left(\frac{|\mathcal{A}^{(g)}|}{|\mathcal{A}_\xi^{(g)}|} \right) = \# \text{ bits pointing out } v \text{ as an exception}$$

Alice



$M =$



Bob



KGIST Anomaly Scores

- Anomalous entities: violate many rules
 - ✧ *MDL intuition*: many bits to describe a node as an exception
- Anomalous triples: unexplained edges ($L(G|M)$) + anomalous endpoints

$$\eta(s, p, o) = \underbrace{\eta(s) + \eta(o)}_{\text{node endpoints}} + \underbrace{\eta^{(p)}(s, p, o)}_{\text{predicate}}$$

Alice



Bob



$$\eta^{(p)}(s, p, o) = \begin{cases} \frac{1}{|\mathcal{A}^-|} * \log\left(\frac{|\mathcal{V}|^2 * |L_{\mathcal{E}}| - |\mathcal{A}_M|}{|\mathcal{A}^-|}\right) & \text{if } \mathcal{A}_{s,o,p}^- = 1 \\ 0 & \text{otherwise} \end{cases} = \# \text{ bits describing unexplained triple}$$

KG_{IST} compresses real KGs significantly

NELL

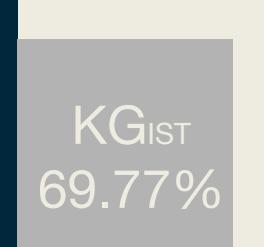


Freq 91.46%



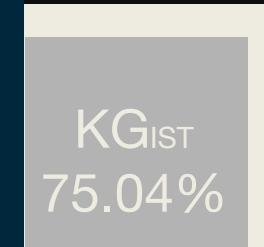
DBpedia

Freq
674.51%



yAGO
select knowledge

Freq
896.33%

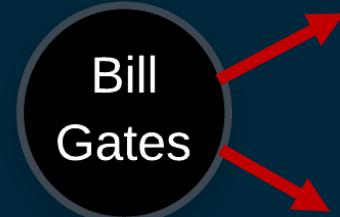


KGIST detects various types of errors

Metric	Supervised		Unsupervised				Select q% of all nodes and, ≤ 0.0188 ≤ 0.0369
	ComplEx	TransE	SDValidate	AMIE+	KGIST_FREQ	KGIST+m	
AUC							
P@100							
R@100							
F1@100							

remove label

billionaire,
entrepreneur,
person



add label

building,
fruit



inject 1 or 2 edges

city

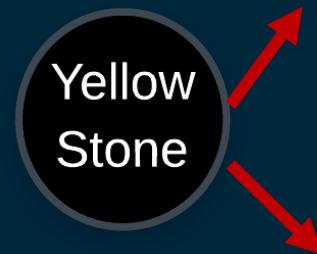
Des Moines



replace label

park,
car

Yellow Stone



KGIST performs best across all types of anomalies.

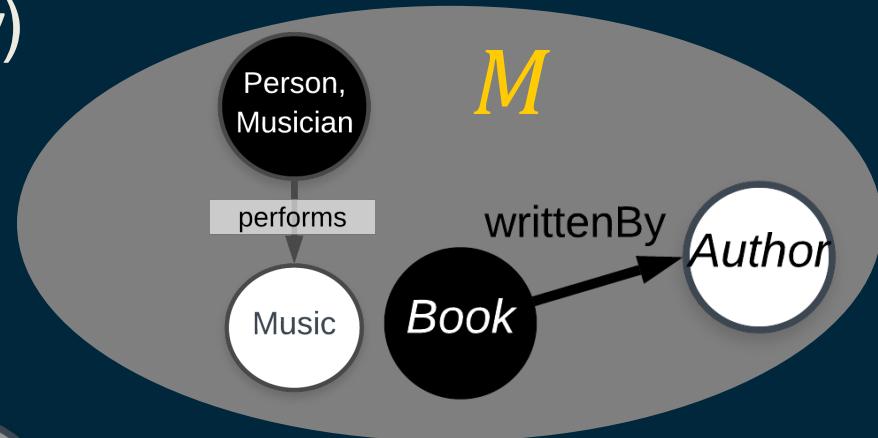
KGIST identifies where information is missing

1. Remove entities / nodes (e.g. *Mary Shelley*)



KGIST identifies where information is missing

1. Remove entities / nodes (e.g. *Mary Shelley*)
2. Run KGIST on perturbed graph
3. Find *where* entities are missing



KGIST identifies where information is missing

Recall for location of missing node

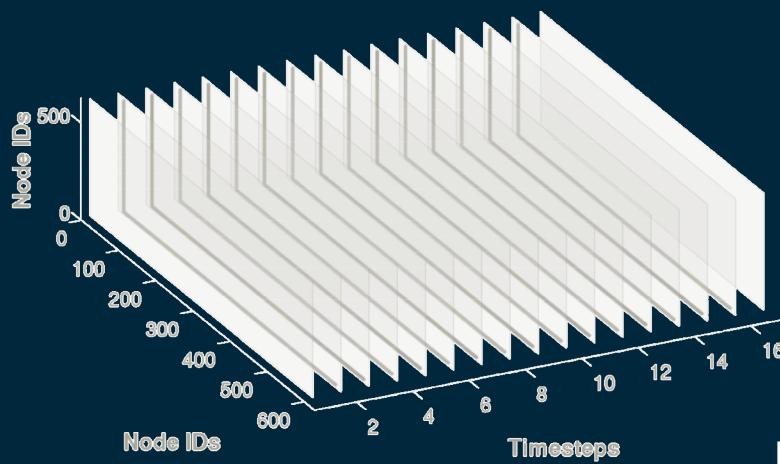
Recall for location + label of missing node

Dataset	Metric	Supervised		Unsupervised	
		LP	AMIE+C [16]	Freq	KGIST
NELL	Recall	N/A	0.6587 ± 0.03	0.4589 ± 0.02	0.7598 ± 0.02
	RL	N/A	N/A	0.3924 ± 0.02	0.6636 ± 0.01
DBpedia	Recall	N/A	0.8187 ± 0.01	0.8049 ± 0.01	0.9288 ± 0.00
	RL	N/A	N/A	0.7839 ± 0.01	0.9179 ± 0.00



KGIST significantly outperforms the baselines. It complements LP methods.

This talk: Summarization Meets Outlier Detection



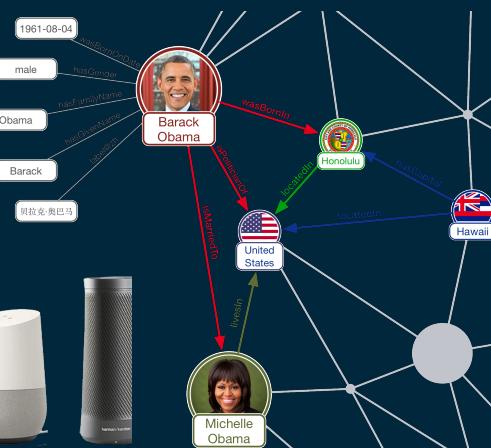
Structural Summaries

[SDM'14, KDD'15,
Dat Bull Eng'17,
SNAM'18,
SDM'19,
KDD'19a, KDD'20...]

Query-on-the-edge + Rule-based

Summaries

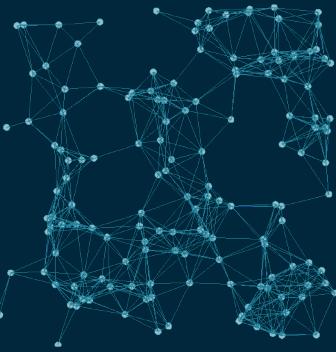
[ICDM'19,
WebConf'20]



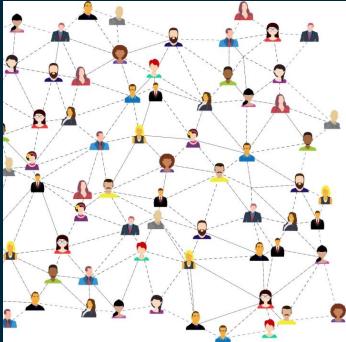
② Graph Streams:
Persistent and bursty activity detection
[KDD'20]

① Knowledge Graphs:
Unified error detection and completion
[WebConf'20]

Summarizing Evolving Networks

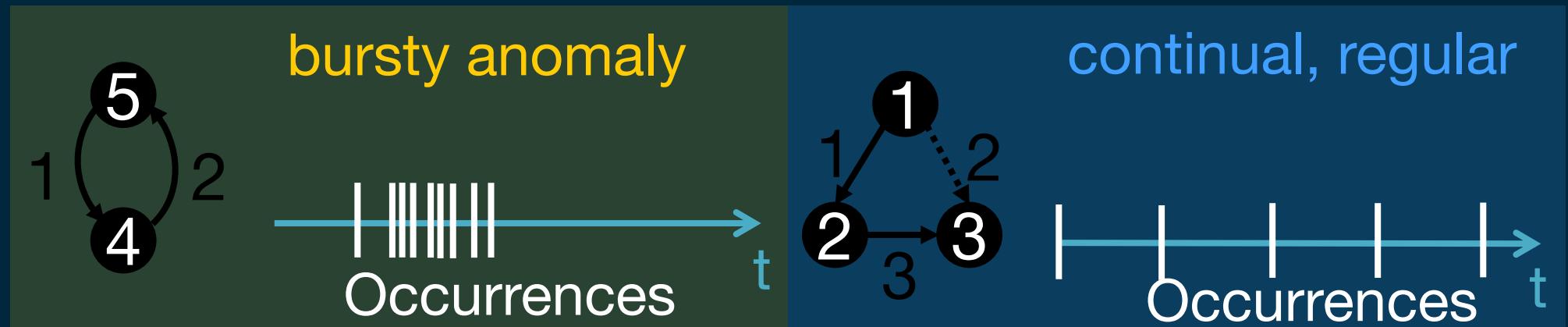
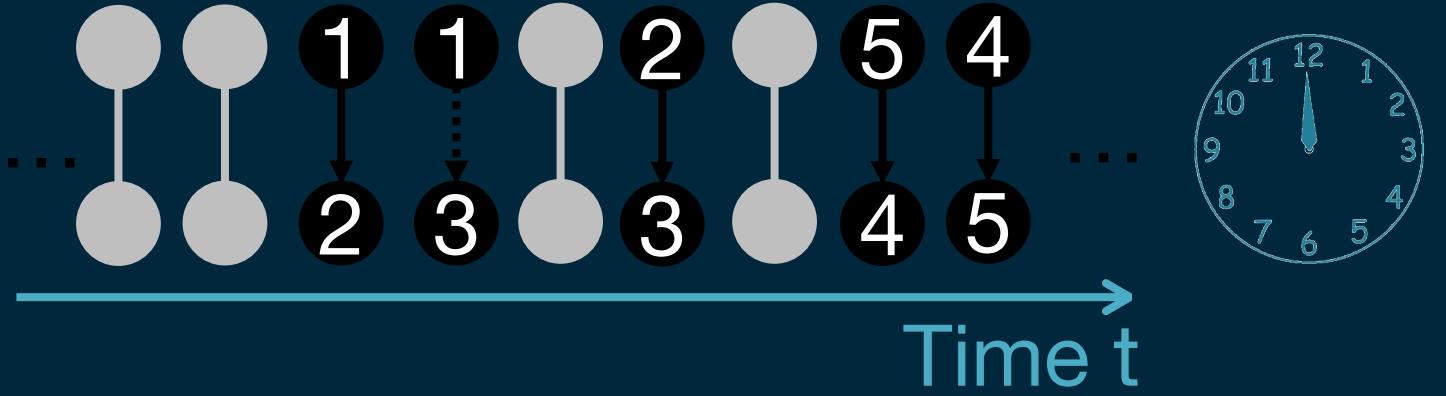
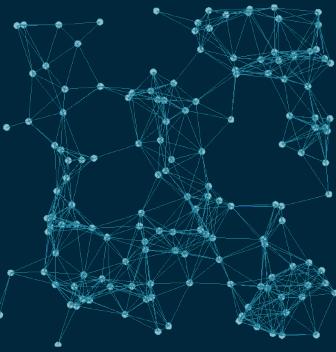


One possibility: summary of frequent graph patterns

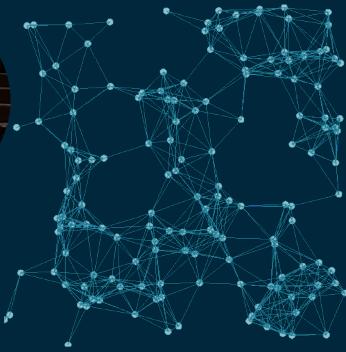


- Related topics:
 - ✧ Motif Mining
 - [Kovanen+, JSTAT'11], [Paranjape+, WSDM'17], [Liu+, WSDM'19]
 - ✧ Frequent Subgraph Mining
 - [Abdelhamid+, TKDE'17], [Aslay+, CIKM'18]

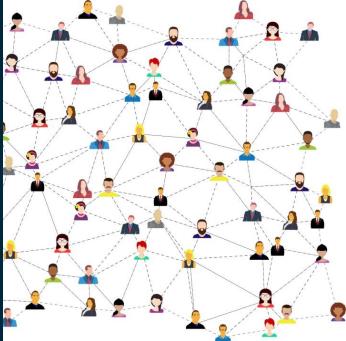
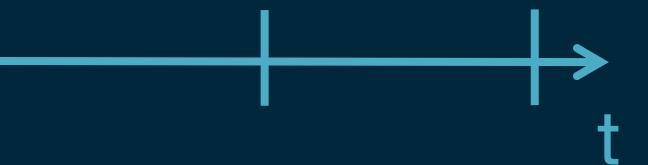
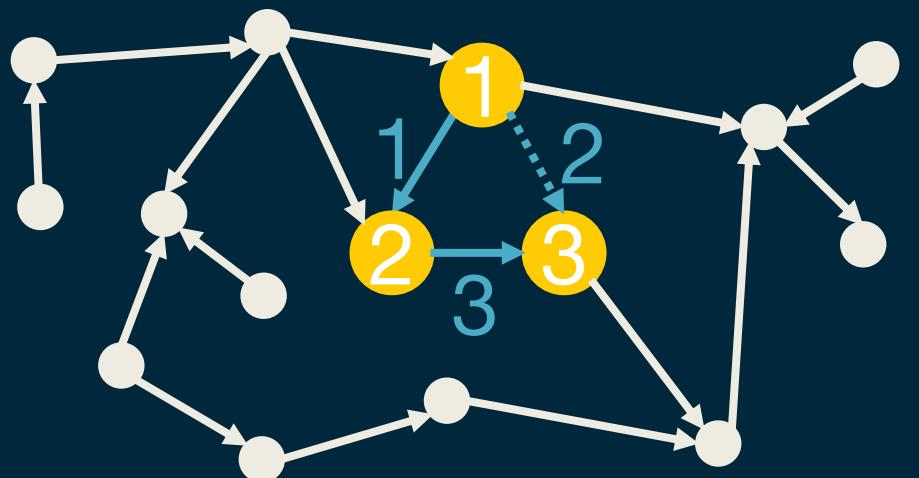
Summarizing Evolving Networks



Summarizing Evolving Networks

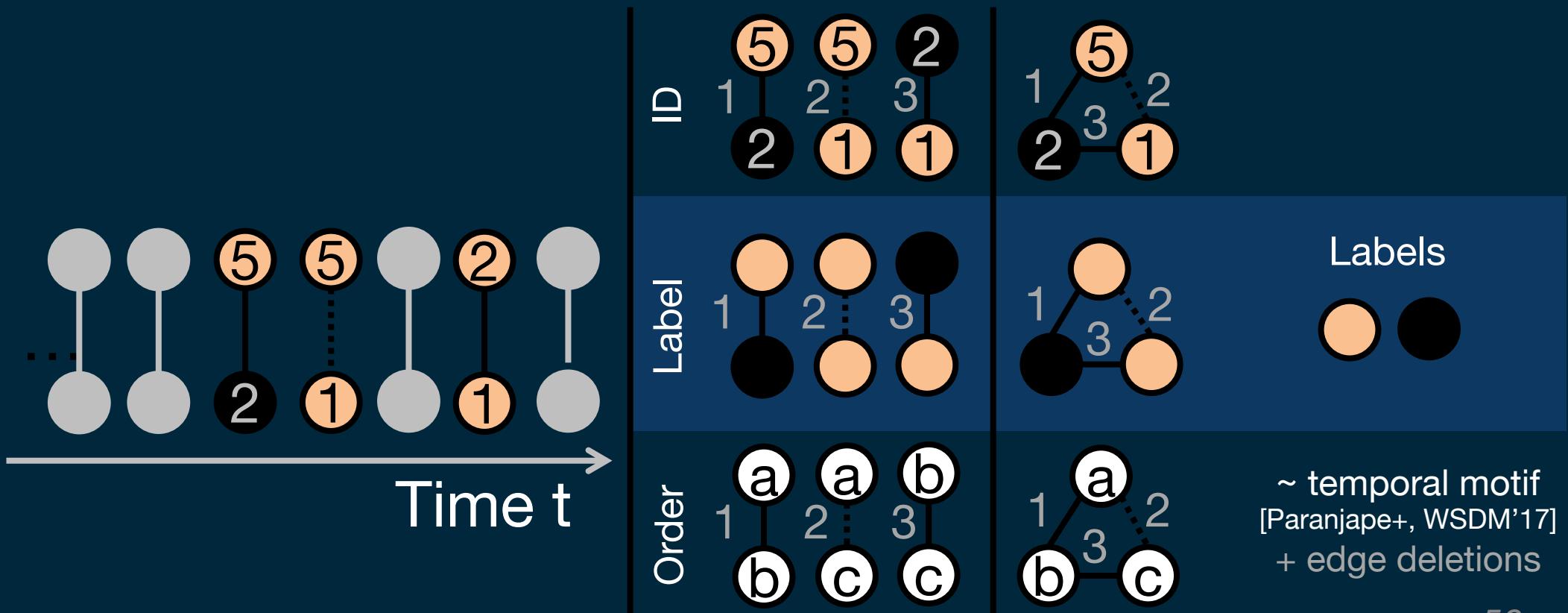


Summarize graph stream G ,
with *persistent* “activity snippets”.



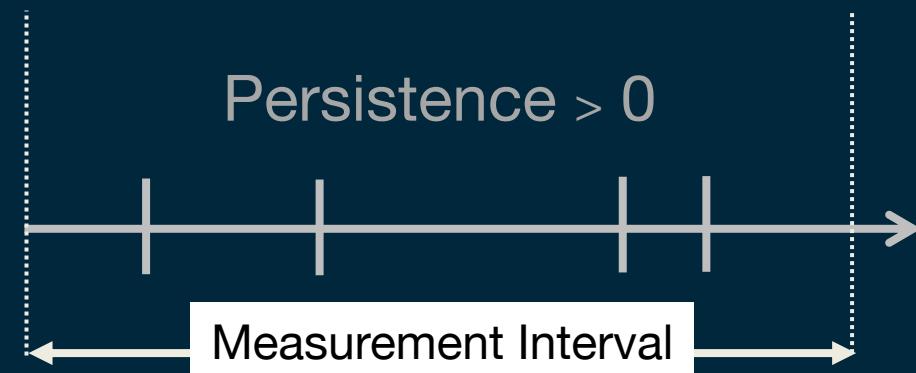
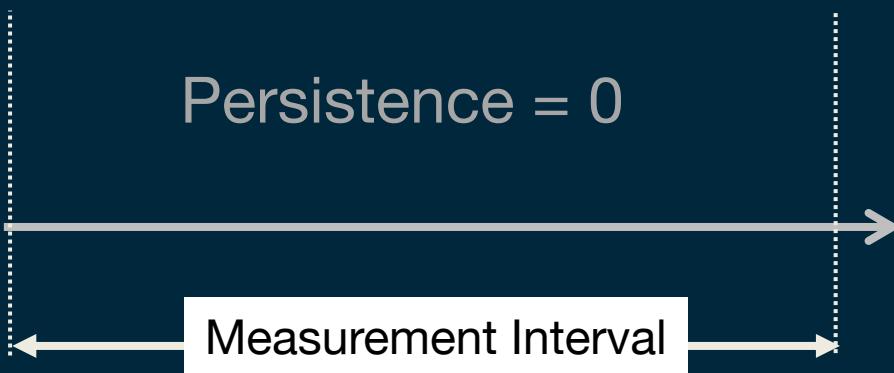
Activity Snippet

An activity snippet describes a sequence of activity among connected nodes in a network



Axioms of Persistence

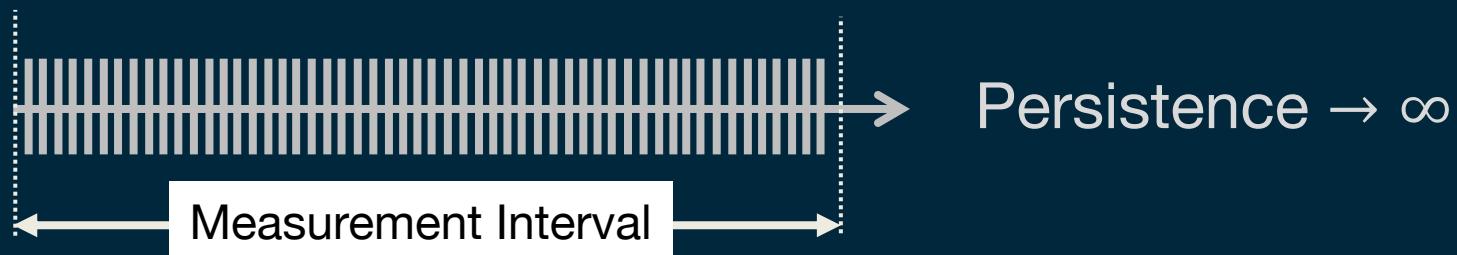
A1: Persistence should be 0 *iff* there are 0 occurrences



Axioms of Persistence

A1: Persistence should be 0 *iff* there are 0 occurrences

A2: As the interval becomes infinitely filled with unique occurrences, persistence should tend to infinity



Axioms of Persistence

A1: Persistence should be 0 *iff* there are 0 occurrences

A2: As the interval becomes infinitely filled with unique occurrences, persistence should tend to infinity

A3: Shifting all occurrences should not affect persistence



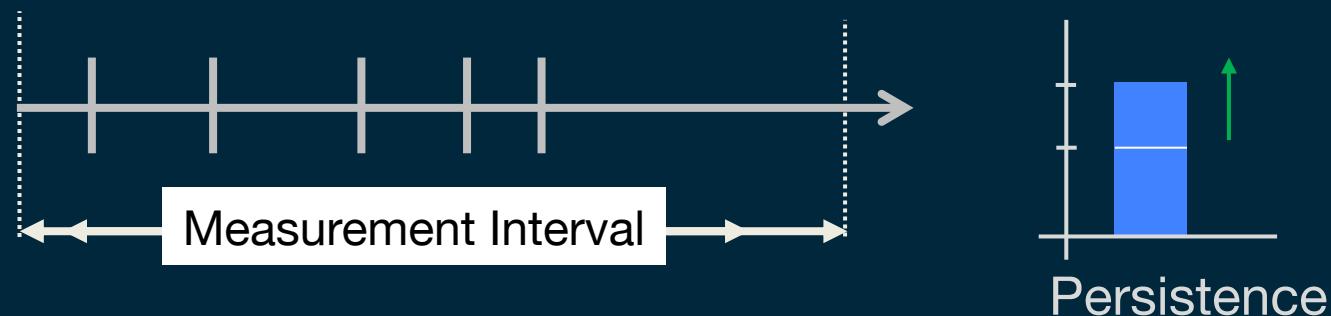
Axioms of Persistence

A1: Persistence should be 0 *iff* there are 0 occurrences

A2: As the interval becomes infinitely filled with unique occurrences, persistence should tend to infinity

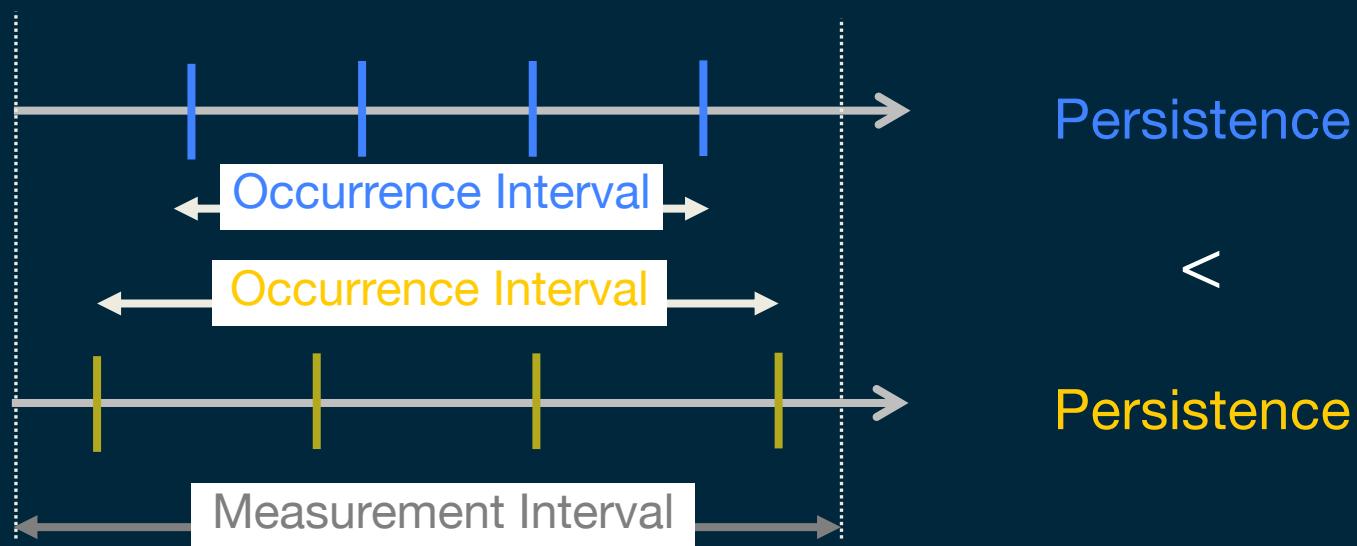
A3: Shifting all occurrences should not affect persistence

A4: Shrinking the interval of measurement leads to higher persistence



Properties of Persistence

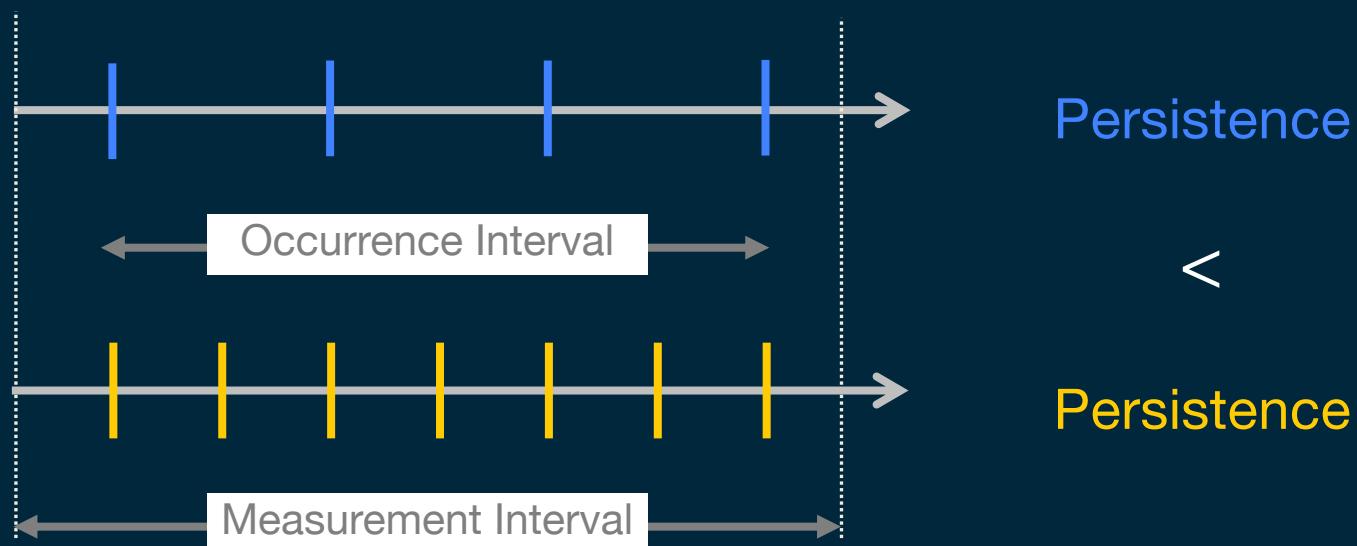
P1: For two snippets with n unique, uniformly-spaced occurrences, persistence is larger for the snippet with occurrences over a wider interval



Properties of Persistence

P1: For two snippets with n unique, uniformly-spaced occurrences, persistence is larger for the snippet with occurrences over a wider interval

P2: For two snippets with unique, uniformly-spaced occurrences spread out over the same interval, persistence is larger for the snippet with more occurrences

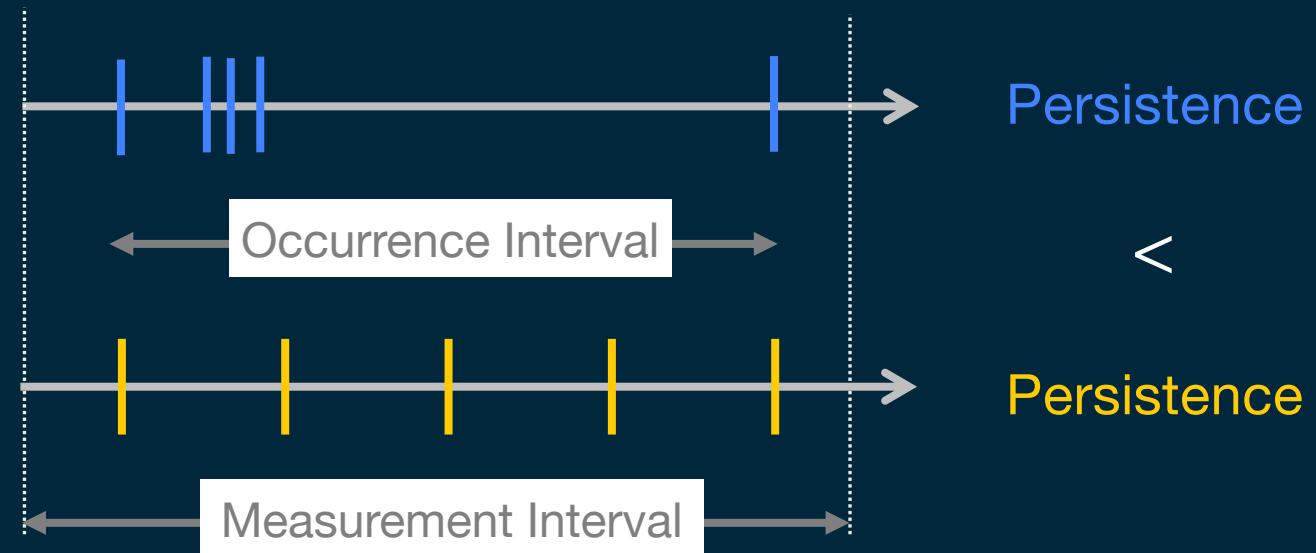


Properties of Persistence

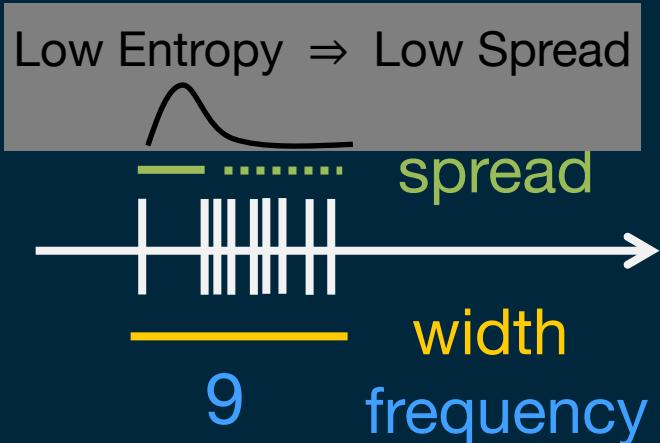
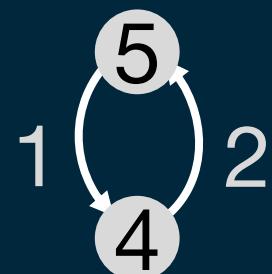
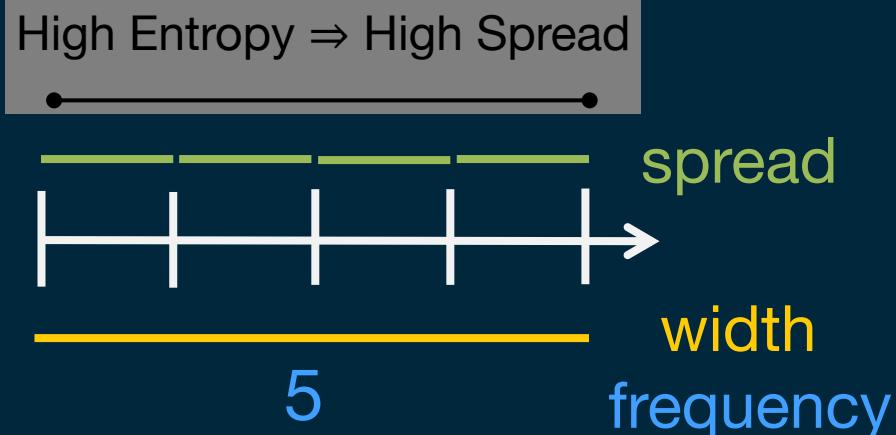
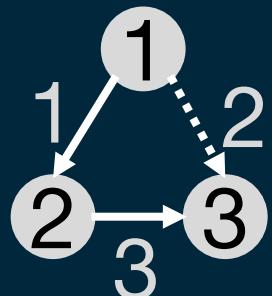
P1: For two snippets with n unique, uniformly-spaced occurrences, persistence is larger for the snippet with occurrences over a wider interval

P2: For two snippets with unique, uniformly-spaced occurrences spread out over the same interval, persistence is larger for the snippet with more occurrences

P3: The persistence of a snippet with n unique occurrences in an interval is maximized *iff* the occurrences are spread out uniformly



Measuring Snippets' Persistence



$$P(x; [t_s, t_e]) \triangleq W(x; [t_s, t_e])^\alpha F(x; [t_s, t_e])^\beta S(x; [t_s, t_e])^\gamma$$

$$\frac{|[t_f, t_l]_x| + 1}{|[t_s, t_e]| + 1}$$

$$\log_{10}(|O_x| + 1)$$

$$\begin{cases} \frac{H(\Gamma_x)}{\log |\Gamma_x|} + 1, & |\Gamma_x| > 1 \\ 1, & |\Gamma_x| \in \{0, 1\} \end{cases}$$

$$H(\Gamma_x) \triangleq - \sum_{g_i \in \Gamma_x} \frac{g_i}{|[t_f, t_l]_x|} \log \frac{g_i}{|[t_f, t_l]_x|}$$

Entropy of distribution
of gaps between
occurrences



Measuring Snippets' Persistence



High Entropy \Rightarrow High Spread

$$P(x; [t_s, t_e]) \triangleq W(x; [t_s, t_e])^\alpha F(x; [t_s, t_e])^\beta S(x; [t_s, t_e])^\gamma$$

Our Solution:
PENminer

Offline + Streaming variants

* Measure of persistence can apply to any
data stream.

width

frequency

9

Low Entropy \Rightarrow Low Spread

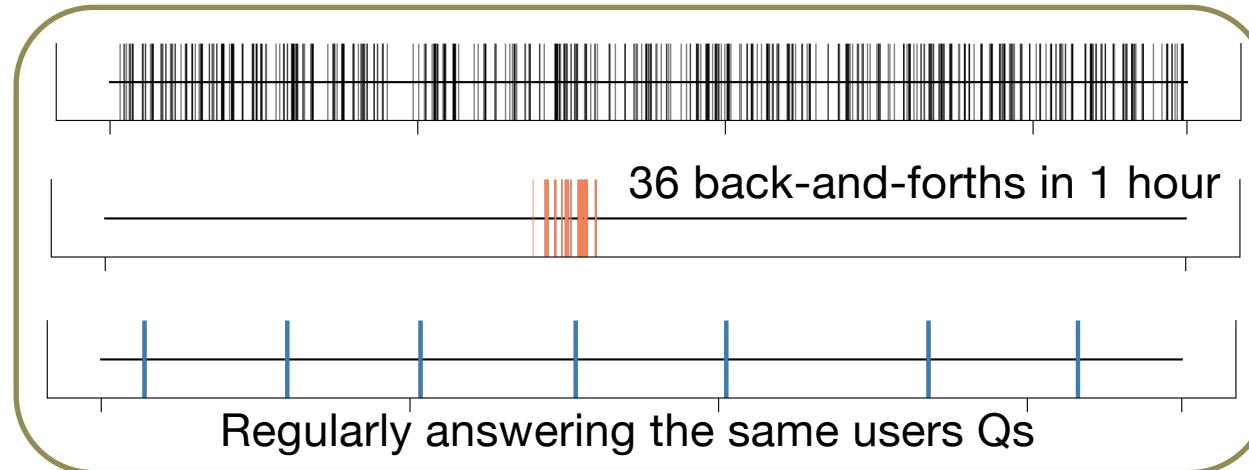
$$(|O_x| + 1)$$

$$\begin{cases} \frac{H(\Gamma_x)}{\log |\Gamma_x|} + 1, & |\Gamma_x| > 1 \\ 1, & |\Gamma_x| \in \{0, 1\} \end{cases}$$

Entropy of distribution
of gaps between
occurrences



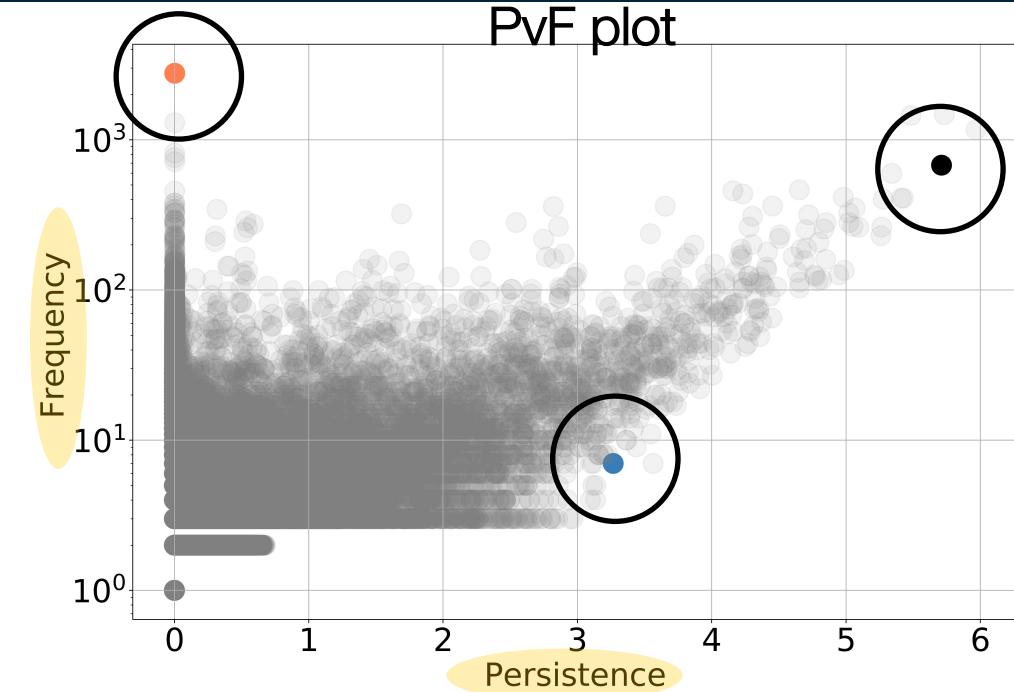
Engaged Discussions & Regular Interactions



u23354
Persistent
(Commented on
Answer)

1 u82199
2 u72603
Bursty
(Commented on
Answer)

u1950
 u55747
Subtly Persistent
(Answered)



Real-time Anomaly Detection

[Eswaran+, ICDM'18]

[Bhatia+, AAAI'20]

[Dai+, PVLDB'14]

Metric	FREQ	SEDANSPOT	MIDAS-R	DS	sPENminer
Subtle	AUC				
	F1@100	subtly anomalous bike trips in Chicago IL			
	F1@1K				
	F1@2K				
Bursty	AUC				
	F1@500K	bursty network attacks in DARPA IP network			
	F1@1M				
	F1@2M				
Avg AUC					

Subtle anomalies like these



Bursty anomalies like these



Setup:

- Represent each activity snippet at time t with a 2d point <frequency, persistence>
- Apply a streaming anomaly detection method (e.g., Random Cut Forest or RCF [Guha+, ICML'16])



Real-time Anomaly Detection

[Eswaran+, ICDM'18]

[Bhatia+, AAAI'20]

[Dai+, PVLDB'14]

Metric	FREQ	SEDANSPOT	MIDAS-R	DS	sPENminer
Subtle	AUC	0.8325 ± 0.02	0.4519 ± 0.01	0.4520 ± 0.02	0.7435 ± 0.03
	F1@100	0.0505 ± 0.01	0.0001 ± 0.00	0.0000 ± 0.00	0.0076 ± 0.00
	F1@1K	0.1812 ± 0.00	0.0035 ± 0.00	0.0003 ± 0.00	0.0378 ± 0.01
	F1@2K	0.1572 ± 0.01	0.0098 ± 0.00	0.0002 ± 0.00	0.0561 ± 0.01
Bursty	AUC	0.8450 ± 0.00	0.6390 ± 0.00	$0.9434 \pm 0.00*$	0.8632 ± 0.00
	F1@500K	$0.3089 \pm 0.00*$	0.2745 ± 0.00	0.3019 ± 0.00	0.3063 ± 0.00
	F1@1M	$0.5351 \pm 0.00*$	0.4527 ± 0.00	0.5274 ± 0.00	0.5295 ± 0.00
	F1@2M	0.7184 ± 0.00	0.6309 ± 0.00	$0.8378 \pm 0.00*$	0.8066 ± 0.00
Avg AUC		0.8388 (2)	0.5455 (5)	0.6977 (4)	0.8034 (3)
				0.8834 (1)	

Subtle anomalies like these



Bursty anomalies like these

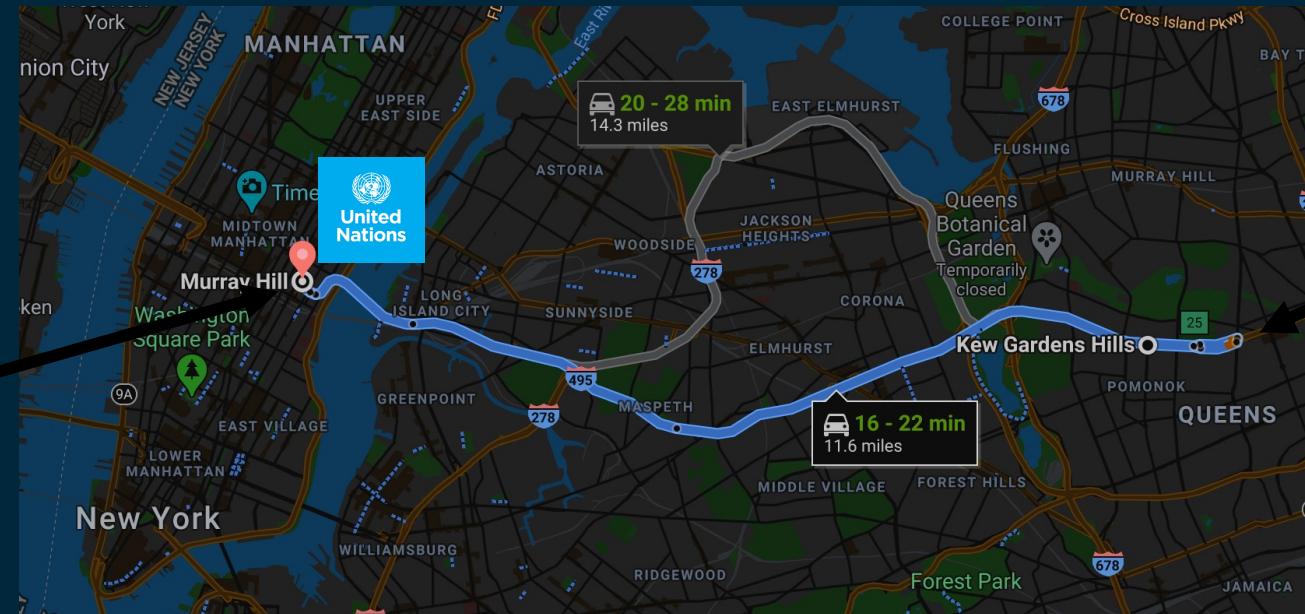


PENminer outperforms all baselines at the new task of finding subtle anomalies, and performs competitively at finding bursty anomalies against baselines designed specifically for that task.

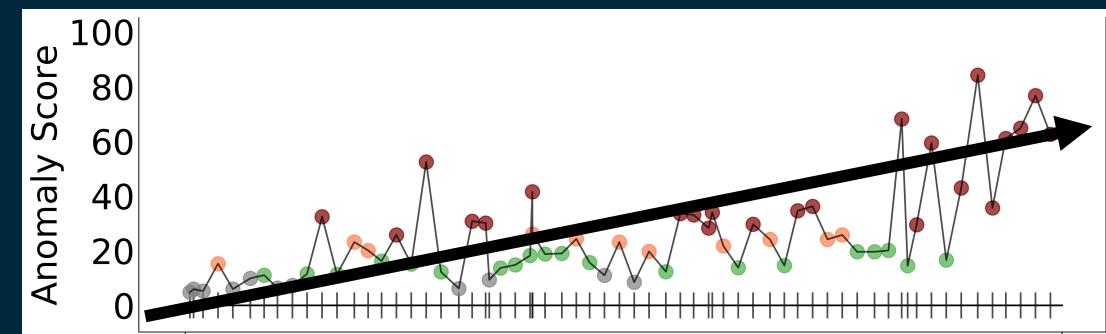


Surprisingly Regular Taxi Trips

Mysterious trip every day
From Queens
To near UN building
Around midnight
For over two months

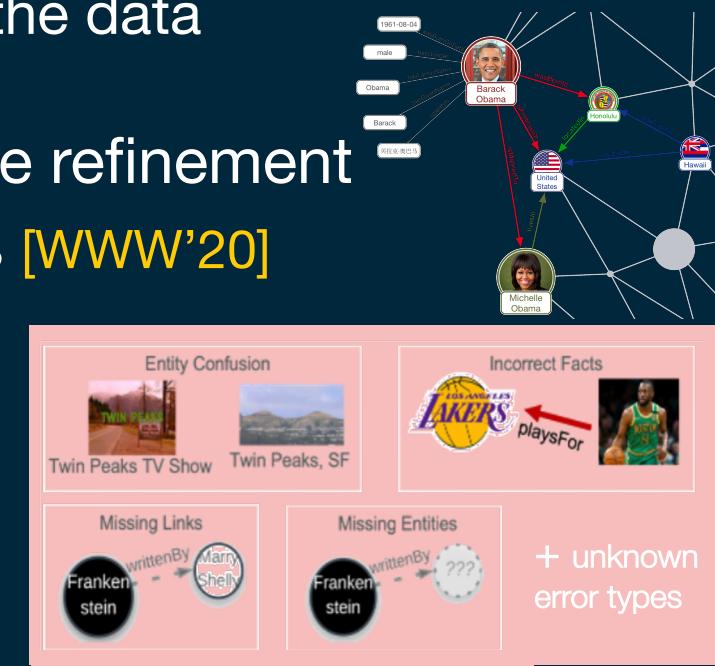


NYC Taxi data



Recap: Graph Summarization Meets Outlier Detection

- Summarization can help identify *patterns* and *anomalies* in the data
- Rule-based summarization of KGs can help unify multiple refinement tasks that are traditionally solved by tailored approaches [WWW'20]
 - ✧ KGist can identify various types of errors in KGs and missing information
- Summarization of graph streams with persistent activity snippets [KDD'20]
 - ✧ Beyond just frequency; capture *how* patterns evolve
 - ✧ The relationship of frequency and persistence highlight anomalies



Talk based on the following papers

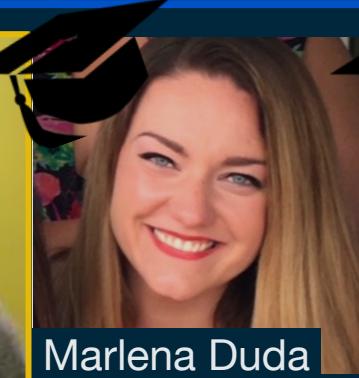
- Y. Liu, T. Safavi, A. Dighe, D. Koutra. [Graph Summarization Methods and Applications: A Survey](#). ACM Computing Surveys 2018.
- Caleb Belth, Xinyi (Carol) Zheng, Jilles Vreeken, Danai Koutra. [What is normal, What is Strange, and What is Missing in a Knowledge Graph: Unified Characterization via Inductive Summarization](#). The Web Conference (WWW/TWC) '20.
- Caleb Belth, Xinyi (Carol) Zheng, Danai Koutra. [Mining Persistent Activity in Continually Evolving Networks](#). ACM SIGKDD '20.

Thank you!
Questions?

Danai Koutra
dkoutra@umich.edu



Caleb Belth



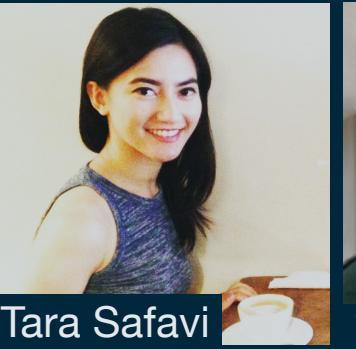
Marlena Duda



Mark Heimann



Di Jin

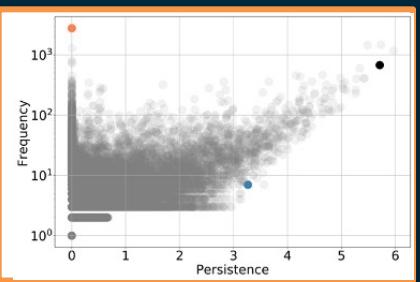
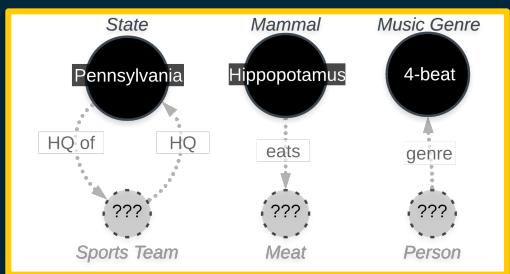


Tara Safavi



Yujun Yan

Graph Summarization Meets Outlier Detection

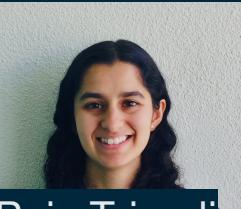


<https://github.com/GemsLab/KGIST>

<https://github.com/GemsLab/PENminer>



Jing Zhu



Puja Trivedi



Jiong Zhu

