Generations of Knowledge Graphs: The Crazy Ideas and The Business

Xin Luna Dong

1/2023

This talk does not represent the company's point of view

Theme I. Three Generations of Knowledge Graphs 3. Media-Rich **1. Entity-Based KGs** 2. Text-Rich KGs

-72



Time-Based KGs



Google Generic KGs

Amazon Product KGs

Theme II. The Recipe from Innovation to Practice

1,000,000,000,000

Feasibility

Quality

A prototype, an experiment, to show it is possible Production quality, E2E MVP experiences

Repeatability

E2E pipelines, automations to allow for extensions to more domains **Scalability**

Significant cost reduction for scale ups

Ubiquity

High coverage, long-tail use cases, assumption removal, to next cycle of inventions

From Roofshots to Moonshots

1,000,000,000,0000

Feasibility

Quality

Repeatability

Scalability

Ubiquity





Roofshots



Generation #1: Entity-Based Knowledge Graphs

Entity-Based KG Example



Entity-Based KGs

Characteristics of Entity-Based KGs

- Ontology (types, relationships) manually defined w. clear semantics
- Entities are named-entities, w. no overlap

Crazy Idea

Create a graph of entities and relationships to represent the world



Transforming Wikipedia to A Knowledge Graph

Feasibility

A prototype, an experiment, to show it is possible



First pot of gold: Transform Wikipedia Infoboxes to knowledge entities and relationships



Transforming Wikipedia to A Knowledge Graph

1 Quality Production

quality, E2E MVP experiences High quality of Wikipedia data guarantees production quality. Common practice in industry

Well-known Examples:



Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weikum. Yago: A core of semantic knowledge. WWW, 2007.

Integrating Data from Different Sources

1M

Repeatability

E2E pipelines, automations to allow for extensions to more domains

Tools for more data sources, long-tail domains

IMDB



ngwriter and actress
ngwriter and actress
ngwriter and actress
ined
itora y actriz mexicana
dit
itora y _{dit}

Are they the same person? → Entity Linkage Are "Born" and "date of birth" the same? → Cchema Algonment Why "May 14, 1982" vs "7 November 1983"? → Data fusion Heterogeneity, Heterogeneity, Heterogeneity

WikiData

Integrating Data from Different Sources

Repeatability

E2E pipelines, automations to allow for extensions to more domains **Biggest Challenge: Entity Linkage**

Random forest on attribute-wise similarity

Results between Freebase and IMDb movies

	Precision	Recall			
Movie	99.0%	98.7%			
People	99.3%	99.6%			
1.5M labels					



Integrating Data from Different Sources

Biggest Challenge: Entity Linkage

Apply active learning to minimize #labels



Zhu et al., Collective multi-type entity alignment between knowledge graphs, WebConf 2020. Zhang et al., AutoBlock: A hands-off blocking framework for entity matching, WSDM 2020.

Zhang et al., OpenKI: Integrating open information extraction and knowledge bases with relation inference, NAACL 2019.

Trivedi et al., LinkNBed: Multi-Graph representation learning with entity linkage, ACL 2018.

Extracting Data from Semi-Structured Websites



Scalability

Significant cost reduction for scale ups





卧虎藏龙 臥虎藏龍 (2000)

导演:李安 编则: 王蕙玲 / 詹姆斯·夏莫斯 / 蔡国荣 主演:周润发/杨紫琼/章子怡/张震/郑佩佩/ 軍多... 类型: 剧情 / 动作 / 爱情 / 武侠 / 古装 制片国家/地区:台湾 / 香港 / 美国 / 中国大陆 语言: 汉语普通话 上映日期: 2000-10-13(中国大陆) / 2000-05-16 (戛纳电影节) / 2000-07-07(台湾) / 2000-07-13 (香港)/2001-01-12(美国) 片长: 120 分钟 又名: Crouching Tiger, Hidden Dragon IMDb链接: tt0190332

评价: 公公公公公 想看

♀ 写短评 ℓ 写影评 + 提问题 分享到 ▼

卧虎藏龙的剧情简介·····

一代大侠李慕白(周润发饰)有退出江湖之意,托付红颜知己俞秀莲(杨紫琼饰)将青冥剑转交给贝勒爷 (郎雄饰)收藏,不料当夜遭玉娇龙(章子怡)窃取。俞秀莲暗中查访也大约知道是玉府小姐玉蛟龙所为,她想 办法迫使玉蛟龙归还宝剑、免伤和气。但李慕白发现了害死师傅的碧眼狐狸(郑佩佩饰)的踪迹、她隐匿于玉府 并收玉蛟龙为弟子。而玉蛟龙欲以青冥剑来斩断阻碍罗小虎(张震饰)的枷锁,他们私定终身。关系变得错综复 杂、俞秀莲和李慕白爱惜玉蛟龙人才难得、苦心引导、但玉蛟龙却使性任气不听劝阻...... ©豆瓣

Two-stage extraction based on distant supervision

- Identify subject
- Identify (attribute, value) pairs

Lockard et al., Ceres: Distantly supervised relation extraction from the semi-structured web. VLDB, 2018.

豆瓣评分

7.9

2星 23%

1星 0.4%

好于 92% 武侠片

好于 90% 动作片

26.7% 45.0%

25.6%

Extracting Data from Semi-Structured Websites



Web Knowledge Extractions & Fusion



Solution: Extract knowledge from different types of web sources, and apply knowledge fusion to remove noises and generate probabilistic knowledge facts

Dong et al., Knowledge Vault: A Web-scale approach to probabilistic knowledge fusion. SIGKDD2014.

Web Knowledge Extractions & Integration

1T

Ubiquity

High coverage, long-tail use cases, assumption removal, to next cycle of inventions

Why NOT in Google Knowledge Graph?

- Didn't reach production quality
 - Accuracy=0.7 is far less than the production requirement (Accuracy=0.99)

#Triples	3.2B (0.3B w. pr>=0.7)
#URLs	2.5B (28M Websites)
#Extractors	16

- Didn't find an **E2E MVP experience**
 - The 0.3B facts (vs. 70B in KG) are very "long-tail" to support meaningful use cases
- But underlying techs applied in long-tail knowledge collection etc.

Dong et al., From data fusion to knowledge fusion. PVLDB2014. Dong et al., Knowledge Vault: A Web-scale approach to probabilistic knowledge fusion. SIGKDD2014. Dong et al., Knowledge-based trust: Estimating the trustworthiness of web sources. PVLDB2015.

Entity-Based KG Summary

Crazy idea: A graph of entities and relationships to represent the world

Main challenges: Heterogeneous data everywhere



What About the Product Domain?

Ubiquity

High coverage, long-tail use cases, assumption removal, to next cycle of inventions

Can we do the same for products?



Generation #2: Text-Rich Knowledge Graphs

Text-Rich KG Example



Text-Rich KGs

Characteristics of Text-Rich KGs

- Ontology (types, relationships) very complex with overlaps and ambiguities; E.g., millions of product types
- Entities may not be named-entities, such as products E.g., "Onus 2 Colors Highlighter Stick, Shimmer Cream Powder Waterproof Light Face Cosmetics, creamy Self Sharpening Crayon STick Highlighter" vs. "Xin Luna Dong"
- Attribute values are oftentimes texts, with overlaps and ambiguities E.g., "Coffee" vs "Cappuccino" as icecream flavors

Crazy Idea:

Finding structure and modeling ambiguity from text sources



Use Case I: Providing Information



Roll over image to zoom in

Brand	Cetaphil
Ingredients	Water, Cetyl Alcohol, Propylene Glycol, Iodopropynyl Butylcarbamate, 2-Bromo-2- Nitropropane-1, 3-Diol, Sodium Lauryl Sulfate, Stearyl Alcohol, Methylparaben, Propylparaben, Sodium Citrate, Butylparaben, Allantoin, Zinc Gluconate.
Scent	Fragrance free
Additional Item Information	Non-Comedogenic, Fragrance- free, Natural
Skin Type	Sensitive

About this item

- Gentle for everyday use; Cetaphil gentle skin cleansing cloths will leave your skin feeling clean, refreshed and balanced after every use
- Removes makeup & dirt: Thoroughly remove makeup and dirt, leaving skin clean
- Mild & non irritating: Soap free formulation won't strip skin of its natural protective oils and emollients



Use Case II: Providing Choices





Use Case III: Improving Search





Use Case III: Improving Search





Use Case III: Improving Search



Use Case IV: Improving Recommendation



KitchenAid KSM	150PSER Artisan Til	t-Head Stand Mixer with	Share 🖂 🚮 🎔
by KitchenAid	-Quart, Empire Red	ad questions	Qty:
List Price: \$429.99 Price: \$249.99 & Fl You Save: \$180.00 (42%)	REE Shipping		\$249.99 + F Only 14 left in st Prei Include 2-Year
i Item is eligible for 6 Mont	Special Financing with your Ama	zon.com Store Card. Learn more	\$14.99
Note: Not eligible for Amaz Amazon.	on Prime. Available with free Prin	ne shipping from other sellers on	Ado
Only 14 left in stock. Estimated Delivery Date:	July 28 - Aug. 2 when you choos	e Standard at checkout.	Turn on 1-Click orde
Color: Empire Red	Premier in easy-to-open package	ng.	Ship to:
			KEVIN DAVENPO
			Add to List
ar But Different			Add to Wedding
More Capacity	More Attachments	Different Brand	Other Seller
			\$264.99
		and a	Sold by: Amazon.com
L.			\$264.99 + Free Shipping Sold by: Marcus AV
2		C.R. C.	\$289.00
See more choices	Sen Sher Options	See Color Options	Sold by: goldentech
	KitchenAid KP26M1XER 6	Hamilton Beach 62022 Felectrics	

Do We Need Different Techniques?

Scott's Cakes Dark Chocolate Toffee Cream Filling Candies with Dark Blue Foils in a 1 Pound Red Roses Box by Scott's Cakes Be the first to review this item





Price: **\$19.95** + \$14.95 shipping You can get 5% back on all Amazon.com purchases with the Amazon Prime Store Card. No annual fee.

Note: Not eligible for Amazon Prime.

In Stock. Ships from and sold by Scott's Cakes.



Get it Thu, Aug 22 - Tue, Aug 27

Qty: 1 V

Get it Tue, Aug 20 - Fri, Aug 23 if you choose paid shipping at checkout.

O Deliver to Yaqing - Seattle 98109

Turn on 1-click ordering

Different challenges: Unstructured and Noisy product data



AutoKnow: Self-Driving Product Knowledge Collection



Dong et al., AutoKnow: Self-driving knowledge collection for products of thousands of types, SigKDD, 2020.



AutoKnow: Self-Driving Product Knowledge Collection



Dong et al., AutoKnow: Self-driving knowledge collection for products of thousands of types, SigKDD, 2020.

Extracting Product Attribute Values



Feasibility

A prototype, an experiment, to show it is possible

name	form	scent
Tide Detergent with Febreze Freshness		
Gain Apple Mango Tango Liquid Laundry Detergent	~	
Gain Joyful Expression Powder Detergent	•	
Tide PODS Original Scent HE Turbo Laundry Detergent Pacs 81-load Tub		
Tide PODS Free & Gentle HE Turbo Laundry Detergent Pacs 35-load Bag		



Generic Knowledge Extraction



Generic KE v.s. Product-Specific KE

Bill Gates founded Microsoft in 1975.

First Aid Beauty Ultra Repair Cream: Vegan and Gluten-Free Intense Moisturizer for Dry Sensitive Skin. Perfect for Skin Conditions and Eczema. Pink Grapefruit (14 ounce) **Solution**: *OpenTag*— Applying deep tagging to identify structured attributes from product titles, descriptions, and bullets

First Aid Beauty Ultra Repair Cream: Vegan and Gluten-Free Intense Moisturizer for Dry Sensitive Skin. Perfect for Skin Conditions and Eczema. Pink Grapefruit (14 ounce) About this Item

ITRA REPAIR CREA

 HIAD TO-TOE: Head-to-tee molistraiser that provides instant self-and long-term hydration for ding. distances and history of the beautiful, whiped texture is instantly aborded with no pravy after-feel. Gragefrait has a bright citrus furti scent that is freeh, juicy and sparking. CLINICALLY PROVEN: Formulated with Colloidal Outward, Shea Burter, Carrande 3 and the Flef and Antioodant Booster, It provides immediate relief and clinically proven to increase hydration by 169% immediately upon application.

Product description
Bannin dirgk kinn with First Ald Beauty's Ultra Repair Cream. Suitable for all skin types, especially drug flags
skin, this hydration wonder leaves skin feeling smooth, hydrated and comfortable after just a single u.e.
Mentioned Attributes: Brand SkinType Scent Quantity

Objects are often not entities

Subject is given

Zheng et al., OpenTag: Open attribute value extraction from product profiles, KDD 2018.



OpenTag Extraction from Product Profiles





OpenTag Extraction from Product Profiles



Captures correlations between **BIOE** tags

Attention

Identifies important terms leading to attribute values

sequence info

Word Embedding

Captures semantics of each token

Zheng et al., OpenTag: Open attribute value extraction from product profiles, KDD 2018.



OpenTag Extraction from Product Profiles



Zheng et al., OpenTag: Open attribute value extraction from product profiles, KDD 2018.


Quality goal: 90% accuracy for product attributes



Repeatability

E2E pipelines, automations to allow for extensions to more domains



Understand domain and attributes, and generate LOTS OF training data Train and fine-tune models



İ

Postprocess extraction results to further improve data quality



Pre-publish evaluation as gatekeeper to guarantee high quality data













OF CANDY

Grocery & Gourmet Food > Candy & Chocolate > Mints

Ø

Price: \$84.99 (\$0.89 / Ounce) + \$16.92 shipping

Pay \$14.17/month for 6 months, interest-free with your Amazon Prime Rewards Visa Card

Flavor Name: Pink

Blue	Green	Orange	Pastel Assort	ment Pink	Red
White	Yellow				
Size: 6 Po	und				
1 Pound	2 Pou	und 3 P	ound 4 Pour	nd 5 Pound	6 Pound

7 Pound 8 Pound 9 Pound 10 Pound

- Love of Candy's huge selection of bulk candy now includes Premium Mint Chocolate Lentils in a variety of bold & striking colors. Available in small to large sizes ranging from 1 to 10 lb bags. These beautiful chocolate morsels feature gourmet, dairy free dark chocolate coated in a crispy and crunchy mint candy shell. Similar to M&M's, these mint chocolate candy lentils are fun, bitesized snacks that can be enjoyed during any occasion.
- Sourced from the most esteemed candy makers from around the world, we've
 put together an extremely broad collection of wholesale candy to fulfill your
 every need. Whether you're in need of candy for vending machines, piñatas or
 candy buffets, you can trust that Love of Candy's got you covered. Our
 consistent product quality and unmatched customer satisfaction have quickly
 made Love of Candy the market's most trusted source of high quality,
 wholesale bulk candy.













Category as input for model training



□Identify **1.77MM** incorrect values for Flavor and Scent for Consumables with **90% precision**

Product	Attr	Value
Love of Candy Bulk Candy - Pink Mint Chocolate Lentils - 6lb Bag	Flavor	Pink
Scott's Cakes Dark Chocolate Fruit & Nut Cream Filling Candies with Burgandy Foils in a 1 Pound Snowflake Box	Flavor	1 lb. snowflake box
Lucky Baby - Baby Blanket Envelope Swaddle Winter Wrap Coral Fleece Newborn Blanket Sleeper Infant Stroller Wrap Toddlers Baby Sleeping Bag (color 1)	Flavor	color 1
ASUTRA Himalayan Sea Salt Body Scrub Exfoliator + Body Brush (Vitamin C), 12 oz Ultra Hydrating, Gentle, Moisturizing All Natural & Organic Jojoba, Sweet Almond, Argan Oils	Scent	vitamin c body scrub - 12oz & body brush
Folgers Simply Smooth Ground Coffee, 2 Count (Medium Roast), 31.1 Ounce	Scent	2Packages (Breakfast Blend, 31.1 oz)

Scaling Up Product Knowledge Extraction



Yan et al., AdaTag: Multi-Attribute Value Extraction from Product Profiles with Adaptive Decoding, ACL 2021.

<u>Option 1. Train a single model?</u> *Train/Test Distribution shift -> Invalid predictions*





<u>Option 1. Train a single model?</u> *Train/Test Distribution shift -> Invalid predictions*



Option 2. Train a model for each category?





Figure 2: Our TXtract architecture for hierarchical multi-task learning.



Figure 2: Our TXtract architecture for hierarchical multi-task learning.

Multi-task Learning



Figure 2: Our TXtract architecture for hierarchical multi-task learning.

Train one model on 4K categories, and improve state-of-the-art by 10.4% in F1, and by 11.7% in coverage

	Title	OpenTag	TXtract
1	Controlled Labs Purple Wraath 90 Servings - Purple Lemonade	flavor: -	flavor: purple lemonade
2	Click - Espresso Protein Drink Vanilla Latte - 16 oz.	flavor: espresso	flavor: vanilla latte
3	Mason Vitamins Melatonin 500 mcg Fast Meltz Tablets, Fruit, 60 Count	flavor: -	flavor: fruit
4	Fashion Glitter Matte Eye Shadow Powder Palette Single Shimmer Eyeshadow	scent: palette	scent: -
5	Baby car seat cover, Nursing covers Breastfeeding cover carseat canopy (Style5)	scent: style5	scent: -

Product Knowl. Extraction from Broader Sources

Ubiquity High coverage, long-tail use cases, assumption removal, to next cycle of inventions shape shape shape Mained Hale SHIMMER STICE SH

Beauty & Personal Care > Makeup > Face > Highlighters & Luminizers

Product knowledge extraction from images, reviews, etc.

Ownest 2 Colors Highlighter Stick, Shimmer Cream Powder Waterproof Light Face Cosmetics, Creamy Self Sharpening Crayon Stick Highlighter

Price: \$9.99 (\$9.99 / Count) FREE Shipping on orders over \$25.00 shipped by Amazon or get Fast, Free Shipping with Amazon Prime & FREE Returns

Get \$60 off instantly: Pay \$0.00 upon approval for the Amazon Rewards Visa Card. No annual fee.

Returnable until Jan 31, 2021 ×

Color: A	
Brand	Ownest
Color	A
Rug Size	Light
tem Form	Stick
kin Tone	Light

About this item

- [Makeup Highlighter] Creamy texture and radiant glow finish, feel beautiful on your skin and they
 gave off the most radiant most amazing highlight ever. These cream highlighter sticks provide a
 subtle glow that is perfect for everyagy use.
- [Clever Design] The top of the stick is flat and wide enough to highlight your checkbones with one swipe, and the edge can be used to apply color precisely to the bridge of the nose, cupid's bow, and brow bone with relative ease.
- [Long Lasting&Waterproof] Long lasting, water-tight, sweatproof. Let you keep the perfect makeup
 all day long without worrying about getting caught in the rain or pushed into a pool.
- [Easy to apply] The shimmer sticks are quick and easy to apply. You just twist up the product and swipe it wherever you would like to illuminate.
- [Perfect Gift for Girts] Two popular everyday colors, suit for various occasions. Unique design, lightweight, easy to carry. Suit for personal use or send as gift.

Lin et al., PAM: Understanding product images in cross product category attribute extraction, KDD 2021.

Product Knowl. Extraction from Broader Sources

1T	Multi-modal kno	wledge	extract	tion most important ro	e le
	Models	P(%)	R(%)	F1(%)	
High coverage, long-tail use cases, assumption removal, to next cycle of inventions	PAM w/o text	79.9	63.4	70.7	
	PAM w/o image	88.7	72.1	79.5	
	PAM w/o OCR	82.0	69.4	75.1	
	РАМ	91.3	75.3	82.5	
	2				

Lin et al., PAM: Understanding product images in cross product category attribute extraction, KDD 2021.

Text-Rich KG Summary

Crazy idea: Finding structure and modeling ambiguity from text-rich sources

Main challenges: Sparse & Noisy structured data everywhere



Unmentioned crazy idea: Automatic taxonomy extraction & construction

Mao et al., Octet: Online catalog taxonomy enrichment with self-supervision. SigKDD, 2020. Zhang et al., Minimally supervised structure-rich text categorization via learning on text-rich networks. Zhang et al., OA-Mine: Open-World attribute mining for e-Commerce products with weak supervision. Webconf, 2022

Generation #3: Media-Rich Time-Based KGs

Background: Virtual Intelligent Assistant

Respond to commands

"Hey Siri, set a timer to 7pm"

"Ok, added to today's reminders"



Background: Virtual Intelligent Assistant

Control devices

"Hey Alexa, turn off bedroom lights"





Background: Virtual Intelligent Assistant

Provide information

"Hey, Google, when did winter start?"

"Winter started on Wed, December 21"



Meta's Assistant

Empowering connection to people and experiences in your life



"Hey Facebook" (double press the button on your controller) "Who's online?"--meet up with friends "Open Boot Scher", jump straight

"Open Beat Saber"--jump straight in the game, and more.



What is An Ideal Virtual Intelligent Assistant?

An *intelligent assistant* should be an agent that **knows you and the world**, can **receive your requests** or **predict your needs**, and provide you the **right services at the right time** with your permission.



What is An Ideal Virtual Intelligent Assistant?

An *intelligent assistant* should be an agent that **knows you and the world**, can **receive your requests** or **predict your needs**, and provide you the **right services at the right time** with your permission.



Evolution of Intelligent Assistant

Chatbot Text input



Voice Asst

Voice input



AR/VR Asst Voice + Visual + Context





What Is Different for An AR/VR Assistant?



From Voice-Only to Multi-Modal

"How tall is Empire State Building?"

"What's the name of this building and how tall is it?"

From Context-Agnostic to Context-Aware

"Remember to buy apples and bananas at the grocery store around the corner"

From Reactive to Proactive (Recommendation!)

"What's the weather today?"

"Today is sunny, 70 degree. Would you like to play your favorite morning music?"

From Server-Side to On-Device

Two Sides of One Coin (1): Great Vehicle for Life Recording

MEMEX (MEMory & EXpansion) by Vannevar Bush (1945)

Two Sides of One Coin (2): Great Vehicle for Personal Assistant Recommendation

Personal KG and Memovoir



Time-Based KGs & Media-Rich Memovoir

Characteristics of a Personal KG & Memovoir

- Rich audio/video, associated w. time, location, etc.
- Modeling
 - entities at an abstract level; e.g., latte art coffee (from different stores), my key for front door
 - activities; e.g., dancing, watching ballet
- Historical–each activity is associated with a timestamp

Crazy Idea

Trace and abstract one's life from rich audios/videos, and use it for life experience Q&A, life journal creation, and life recommendation.

Research Problems for Personal KG and Memovoir

- How to record one's life with hardware (memory, battery) constraints?
- How to extract personal knowledge from the recordings?
- What is the best frequency, granularity, and domain richness to capture one's life?
- How to leverage the Personal KG and Memovoir for utility, memoir, and inspiration applications?
- How to leverage Personal KGs to best understand contexts?
- How to combine public and personal KGs for context-aware recommendation?



Take-Aways I. 3 Generations of KGs

1. Entity-Based KGs

2. Text-Rich KGs

Resolving heterogeneity with entity linkage and web knowledge extraction

Extractions and **cleanings** from sparse and noisy source data, and handling semantics ambiguities 3. Media-Rich Time-Based KGs

Many new challenges for knowledge collection and applications

Take-Away 2. From Roofshots to Moonshots

1,000,000,000,0000

Feasibility

Quality

Repeatability

Scalability

Ubiquity





Roofshots



Shameless Advertisements

Book (2021)

Foundations and Trends[®] in Databases 10:2-4

Machine Knowledge Creation and Curation of Comprehensive Knowledge Bases

Gerhard Weikum, Xin Luna Dong, Simon Razniewski and Fabian Suchanek

Two benchmark datasets

- DI2KG Challenge:
 - http://di2kg.dia.uniroma3.it/ #challenge
- Extended SWDE benchmark: <u>https://homes.cs.washington</u> <u>.edu/lockardc/expanded_sw</u> de.html

Thank You

Q&A?