



# Conceptualization for Short Text Understanding

Zhongyuan Wang (王仲远)

**\*Joint work with Haixun Wang, Jun Yan, Yanghua Xiao, Ji-Rong Wen, and many interns**

微软亚洲研究院 数据挖掘与企业智能组  
Data Mining and Enterprise Intelligence Group  
Microsoft Research Asia | Dec 27, 2015

# Short Text

- Search
- Ad keywords
- Anchor text
- Document Title
- Caption
- Question

Short text is *sparse, noisy,*  
*and ambiguous*

# The big question

- How does the mind get so much out of so little?
- Our minds build rich models of the world and make strong generalizations from input data that is *sparse, noisy, and ambiguous* – in many ways far too limited to support the inferences we make.
- How do we do it?



*Science* **331**, 1279 (2011);

# How to Grow a Mind: Statistics, Structure, and Abstraction

Joshua B. Tenenbaum,<sup>1\*</sup> Charles Kemp,<sup>2</sup> Thomas L. Griffiths,<sup>3</sup> Noah D. Goodman<sup>4</sup>



MIT



CMU



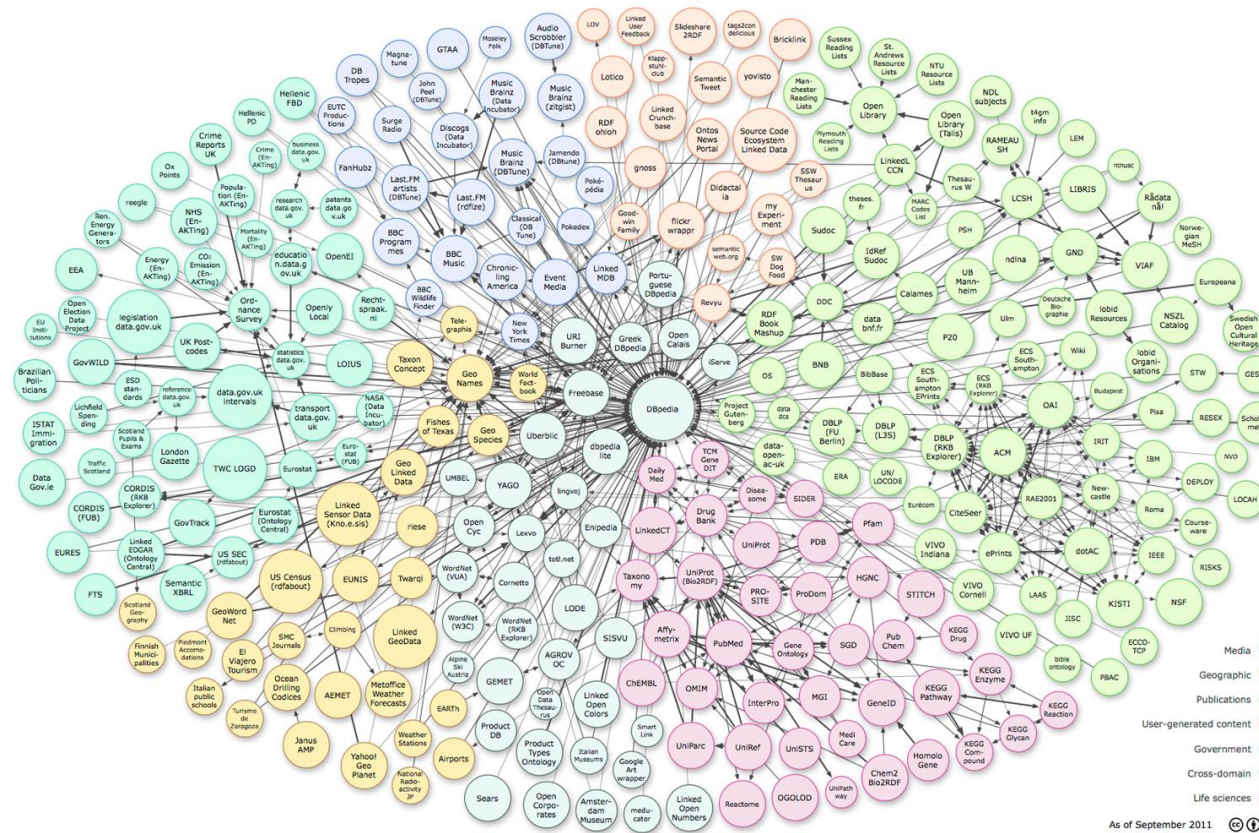
Berkeley



Stanford

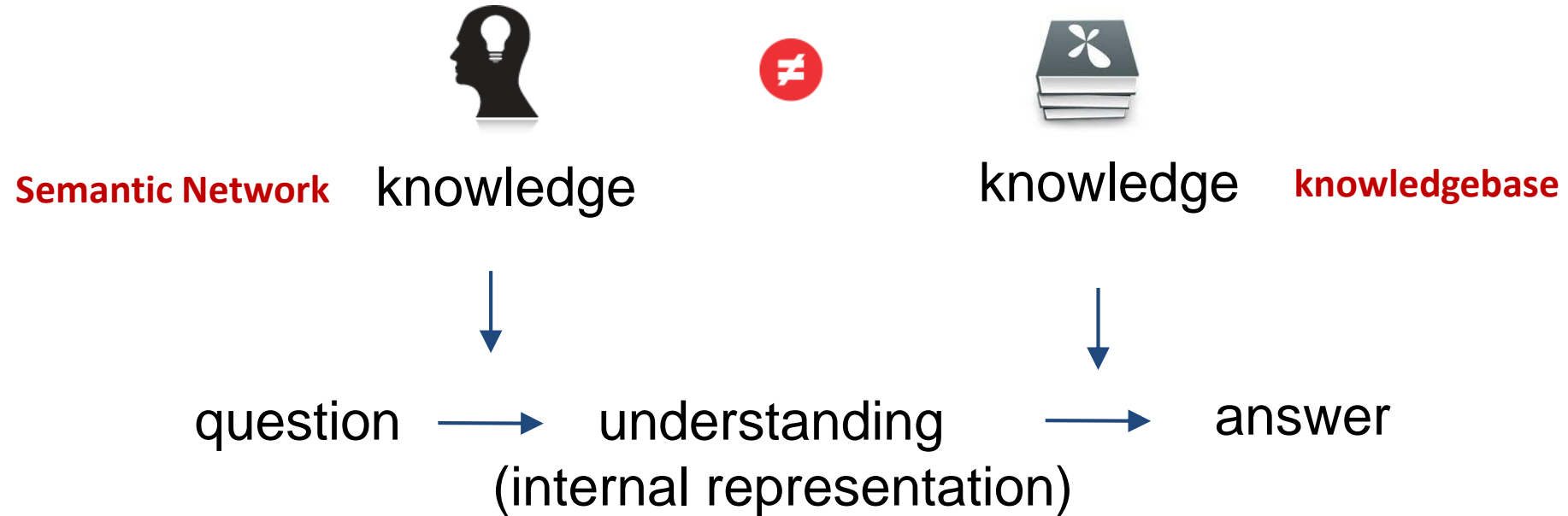
If the mind goes beyond the data given,  
*another source of information* must make up  
the difference.

# Knowledge Base Efforts



[http://lod-cloud.net/versions/2011-09-19/lod-cloud\\_colored.png](http://lod-cloud.net/versions/2011-09-19/lod-cloud_colored.png)

1. *“Python Tutorial”*
2. *“Who was the U.S. President when the Angels won the World Series?”*



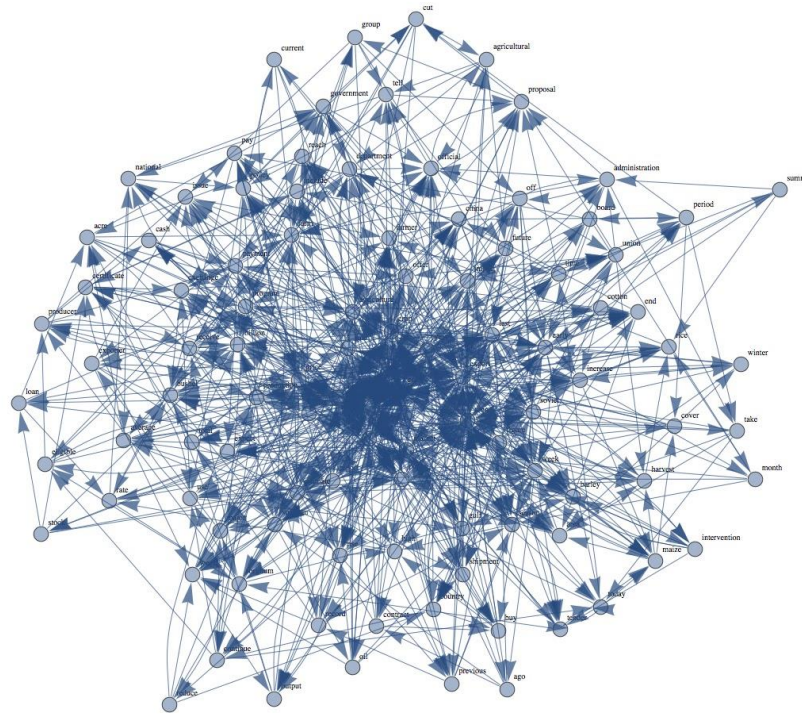
# Semantic Network vs. Knowledgebase

Semantic Network	Knowledgebase
Common/linguistic knowledge	Entities Facts
isA isPropertyOf co-occurrence ...	DayOfBirth LocatedIn SpouseOf ...
Typicality, basic level of categorization	Black or White Precision
KnowItAll, Probase	Freebase, Yago



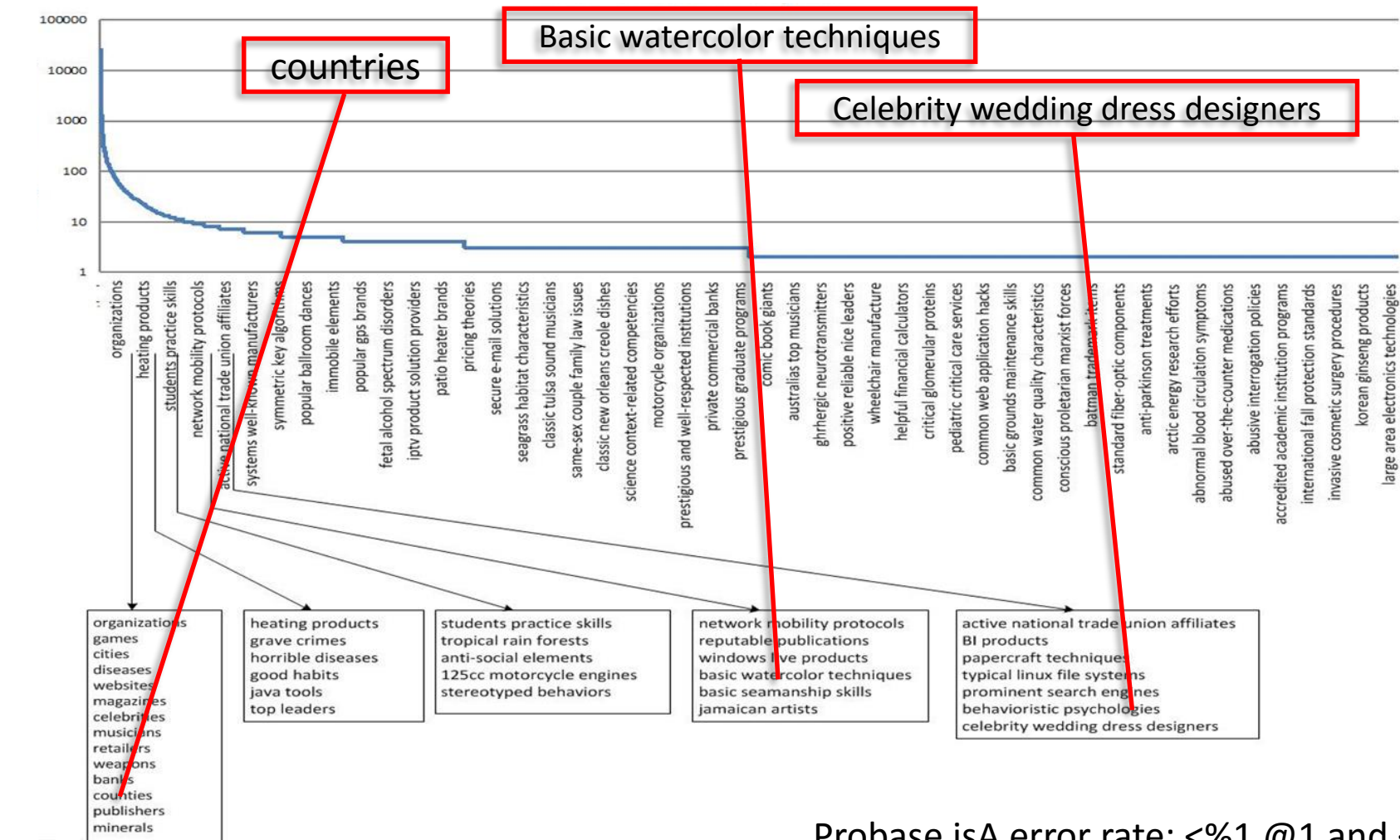
# Probase: A Semantic Network

<http://research.microsoft.com/probase/>



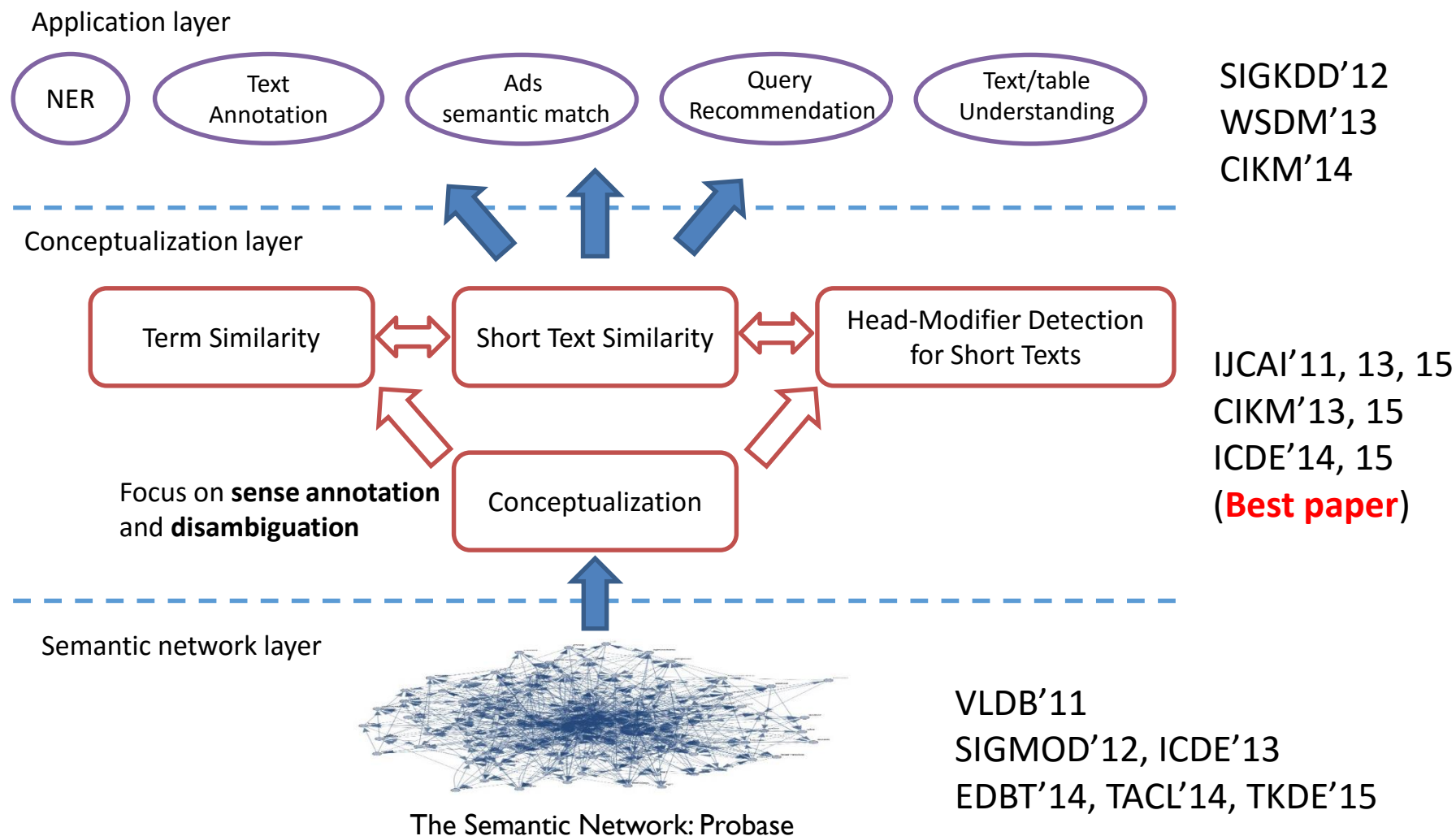
<b>Nodes:</b>	<b>Concepts</b> (“Spanish Artists”)	<b>Entities</b> (“Pablo Picasso”)	<b>Attributes</b> (“Birthday”)	<b>Verbs/Adjectives</b> (“Eat”, “Sweet”)
<b>Edges:</b>	<b>isA</b> (concept, entities)	<b>isPropertyOf</b> (attributes)	<b>Co-occurrence</b> (isCEOof, LocatedIn, etc)	

# Probase Concepts (2.7 million+)



Probase is A error rate:  $<1\%$  @1 and  $<10\%$  for random pair

# Research Roadmap



# What is short text understanding?

# Add Common Sense to Computing

**Pablo Picasso**

**25 Oct 1881**

**Spanish**

China

Brazil

India

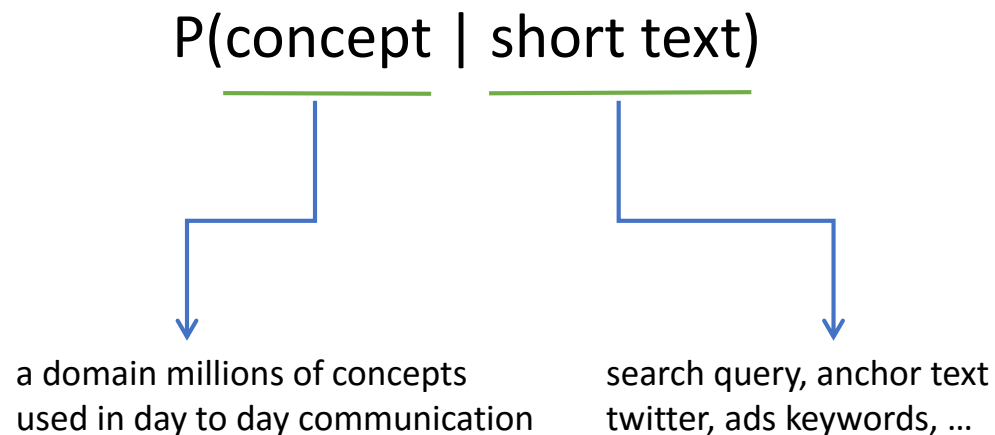
*emerging market*

The engineer is eating an apple

*IT company*

# Conceptualization On Knowledge Engineering (COKE)

- **Conceptualization:** An **explicit** representation for the **short text**



- **Short text** is *sparse, noisy, and ambiguous*
- **Explicit** means
  - Conceptualization results can be *easily understood* by human beings
  - Conceptualization model can be *easily customized* for different scenarios



# Conceptualization On Knowledge Engineering (COKE)

- **Conceptualization:** An **explicit** representation for the **short text**

ShortText: pear apple

[Show Parameters](#)

Elapsed Time = 00:00:00.0140014

pear		apple	
[25/fruit]		[25/fruit]	
25/fruit	0.5724769	25/fruit	0.5718007
fruit	0.0562538	fruit	0.1523192
fresh fruit	0.02918897	fresh fruit	0.05945546
tree fruit	0.01260049	tree fruit	0.01657355
dried fruit	0.01165293	dried fruit	0.01593271
seasonal fruit	0.01160144	seasonal fruit	0.01553914
juice	0.01062348	juice	0.01546597
hard fruit	0.01011614	fruit juice	0.01309836
climacteric fruit	0.009254614	hard fruit	0.01297377
fruit juice	0.009048668	climacteric fruit	0.01062575
sweet fruit	0.008924312	sweet fruit	0.01033749
9405/food	0.1058999	9405/food	0.1241537
food	0.03129783	food	0.06844553
high fiber food	0.008663075	high fiber food	0.01028461
ingredient	0.007349567	ingredient	0.009757149
high-fiber food	0.004699597	fresh food	0.004135756
fresh food	0.004038723	hard food	0.004045118
hard food	0.003555713	high-fiber food	0.003733333
fit	0.00375		
fit	0.00333		
low	0.00344		
low	0.003483		

“pear apple”

ShortText: ipad apple

Conceptualize

[Show Parameters](#)

Elapsed Time = 00:00:00.6670667

ipad		apple	
[15/mobile device/device]		[1/technology company/company]	
15/mobile device/device	0.8072805	1/technology company/company	0.9623328
mobile device	0.01889746	technology company	0.005182603
apple device	0.01723156	computer manufacturer	0.005060604
tablet device	0.01718674	tech company	0.004826283
ios device	0.01706666	innovative company	0.00475833
portable device	0.01549841	computer company	0.004576935
gadget	0.0121459	tech giant	0.004452952
handheld device	0.0105837	technology giant	0.004435823
digital device	0.01037961	successful company	0.00422191
multimedia device	0.009883645	tech stock	0.004145819
wireless device	0.009499655	software company	0.004118673
3/apple product/product	0.1443738	1053/laptop brand/top brand name/brand	0.02506031
apple product	0.01561871	laptop brand	0.0004202249
apple's product	0.005314387	iconic brand	0.0004185871
electronic product	0.004432585	great brand	0.0004146745
hot new apple product	0.004225718	global brand	0.0004121742
apple's high technology product	0.004222373	big brand	0.0004073361
popular apple product	0.004221935	strong brand	0.0003894416
iconic product	0.004196884	popular brand	0.0003739756
revolutionary product			0.0003603094
digital product			0.0003555858
popular product			0.0003552496

“ipad apple”

# Recap: Conceptualization

Conceptualization 1.0  
[IJCAI'11, CIKM'13]:  
*mapping terms to concept  
space based on **Bayesian  
Inference***

Conceptualization 2.5  
[IJCAI'15]: *leveraging  
**verbs, adjective, attribute,**  
etc.*

Conceptualization 2.0  
[ICDE'14, CIKM'14,  
**ICDE'15(Best Paper)**]:  
*incorporating **co-  
occurrence network***

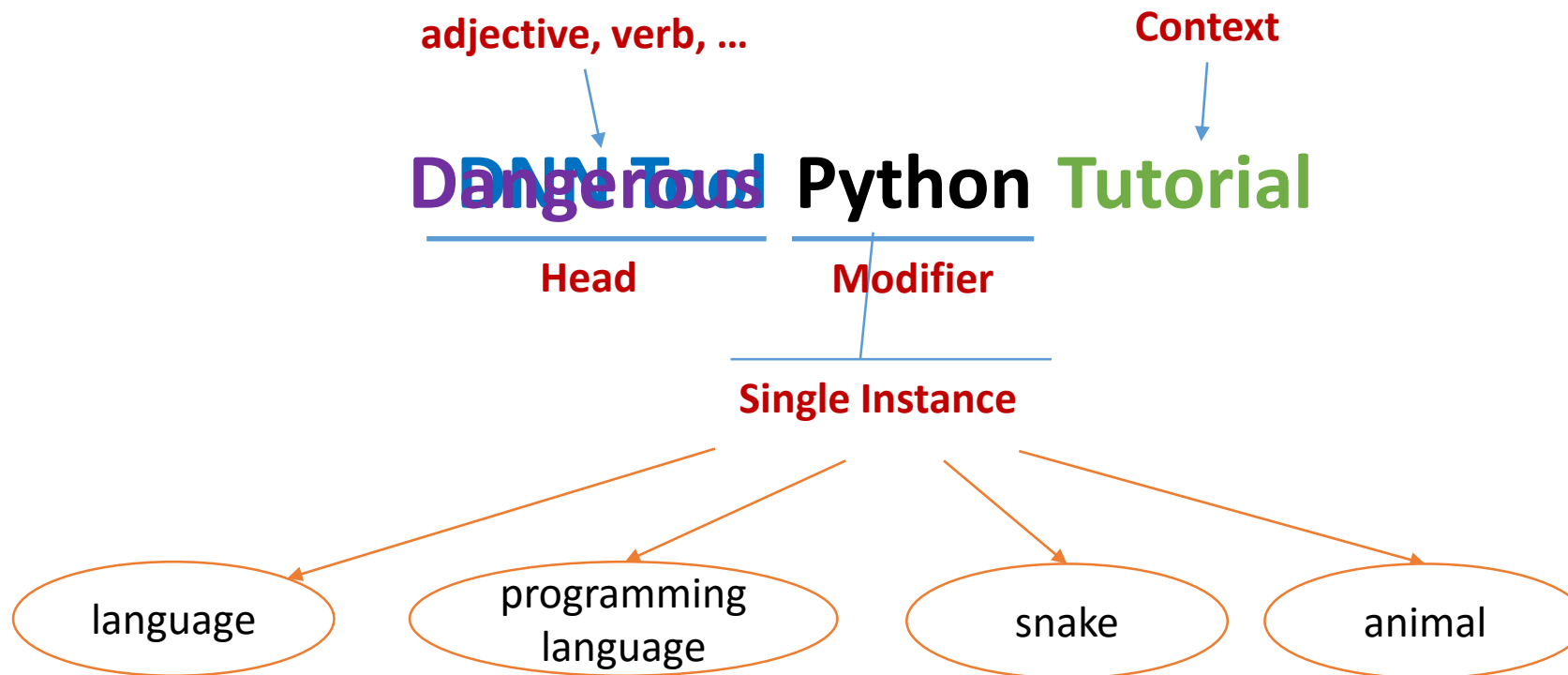
Conceptualization 3.0  
[CIKM'15]: learning-  
based conceptualization/  
leverage **embedding**

## Production Impacts (Shippings):

- Ads relevance (2012)
- MSN Query Recommendation (2012)
- Bing Image Search (2013)
- Table understanding in Power Query (2013)
- Definition Answer in EQnA (2014, 2015)

# What we resolved?

- Short Text Understanding



# Short Text Understanding

- If the short text is a **single instance**...
  - *SIGMOD 2012, CIKM 2015*
- If the short text has **context** for the instance...
  - *IJCAI 2011/2013, ICDE 2015*
- If the short text contains **verb, adjective**...
  - *IJCAI 2015*
- If the short text contains **multiple instance**...
  - *ICDE 2014*
- **Applications**
  - *WSDM 2013, CIKM 2013/2014*

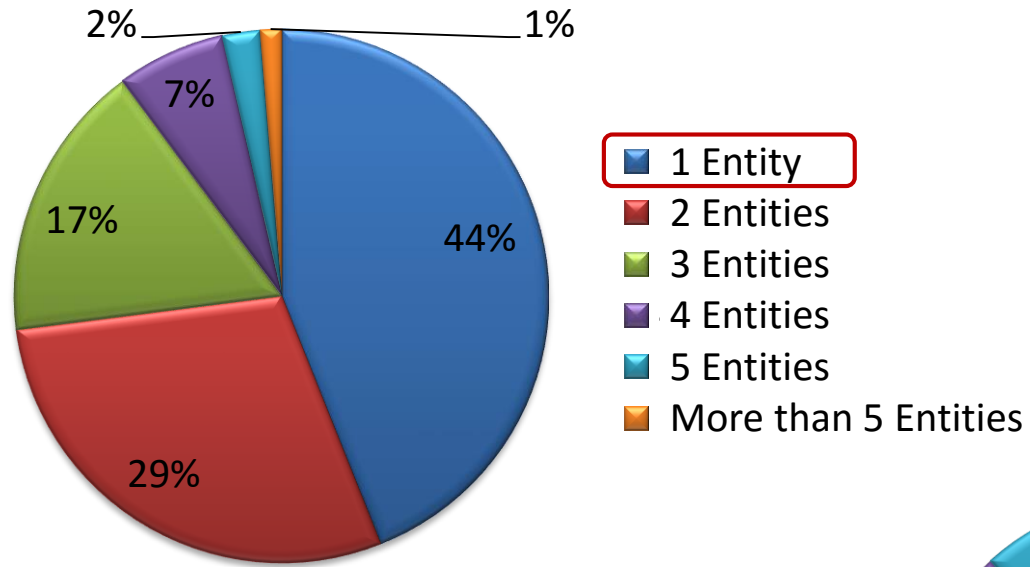


If the short text is a single instance...

“Python”

- Wentao Wu, Hongsong Li, Haixun Wang, and Kenny Zhu, [Probase: A Probabilistic Taxonomy for Text Understanding](#), in *ACM International Conference on Management of Data (SIGMOD)*, May 2012.
- Zhongyuan Wang, Haixun Wang, Ji-Rong Wen, and Yanghua Xiao, [An Inference Approach to Basic Level of Categorization](#), in *ACM International Conference on Information and Knowledge Management (CIKM)*, October 2015.

# Statistics of Search Queries

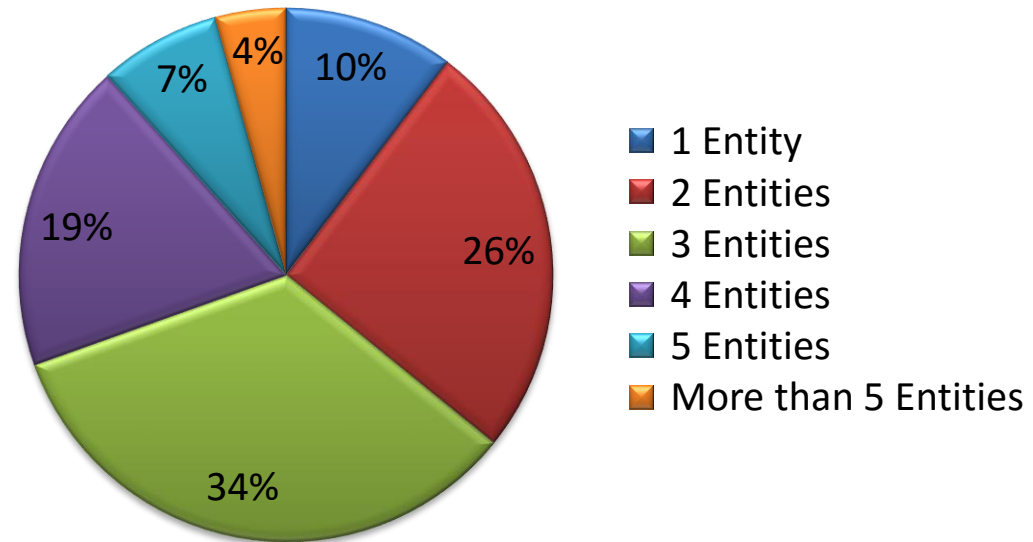


(a) By traffic


angry birds windows phone 8


Entity 1

Entity 2



(b) By # of distinct queries





Web

News

Images

Maps

Shopping

More ▾

Search tools


About 1,310,000,000 results (0.39 seconds)

Microsoft – Official Home Page

[www.microsoft.com/](http://www.microsoft.com/) ▾ Microsoft Corporation ▾

At Microsoft our mission and values are to help people and businesses throughout the world realize their full potential.

Results from microsoft.com



Download Center

Microsoft Download Center: Find the latest downloads for ...

Windows

Downloads - Internet Explorer - Windows 7 - Support - Apps - ...

Support

Microsoft Help and Support provides support for Microsoft ...

Security

Microsoft Safety Scanner - Get security updates - Internet Security


Microsoft Security Essentials

Find out how Microsoft Security Essentials helps guard your PC ...

Surface

Surface Pro 3 - Compare Surface Tablets - At School - Surface RT

In the news



Hands On With Microsoft's Surface 3

TechCrunch - 11 hours ago

Microsoft is back at the well with a new Surface device, the Surface 3. If you're familiar with ...

Microsoft unveils Lumia 640 and Lumia 640 XL for India


GSMarena.com - 1 hour ago

Microsoft launches program to hire people with autism

CNET - 8 hours ago

[More news for Microsoft](#)

Microsoft Corporation




Computer software company

Microsoft Corporation is an American multinational corporation headquartered in Redmond, Washington, that develops, manufactures, licenses, supports and sells computer software, consumer electronics and personal computers and services. [Wikipedia](#)

Stock price: **MSFT** (NASDAQ) \$41.55 +1.25 (+3.11%)  
Apr 6, 4:00 PM EDT - Disclaimer

CEO: [Satya Nadella](#)


Founded: April 4, 1975, [Albuquerque, NM](#)


Customer service: 1 (800) 642-7676 


Headquarters: [Redmond, WA](#)


Founders: [Bill Gates](#), [Paul Allen](#)


Profiles

 Twitter


 Facebook


 LinkedIn


 Instagram


 Google+


People also search for

 Nokia

 Logitech

 Dell


 Sony Corporati...

 Apple Inc.

[View 15+ more](#)

[Feedback](#)

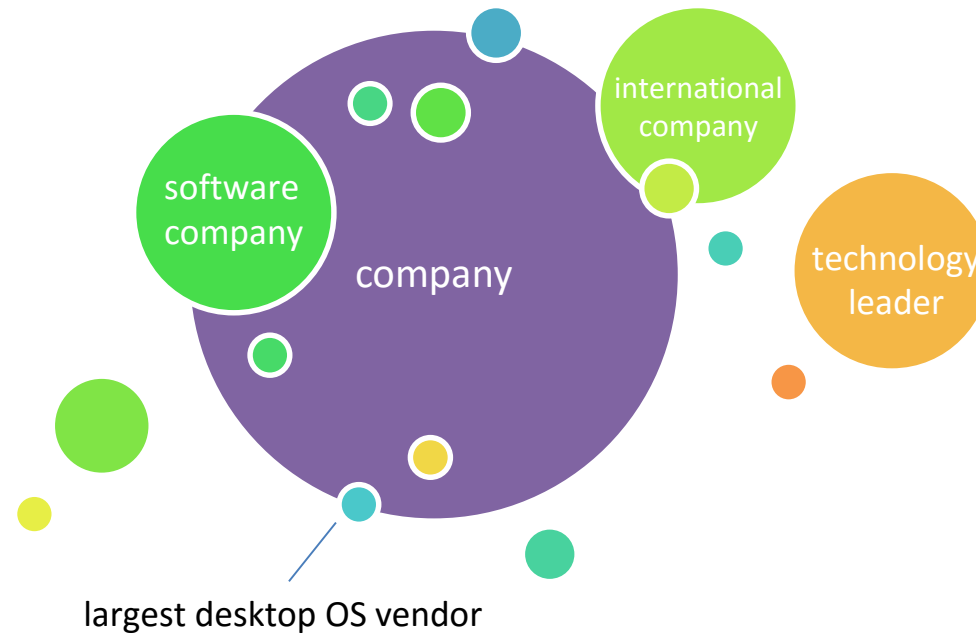
## Knowledge Panel



Microsoft

微软亚洲研究院 数据挖掘与企业智能组 Data Mining and Enterprise Intelligence Group

# A Concept View of “Microsoft”

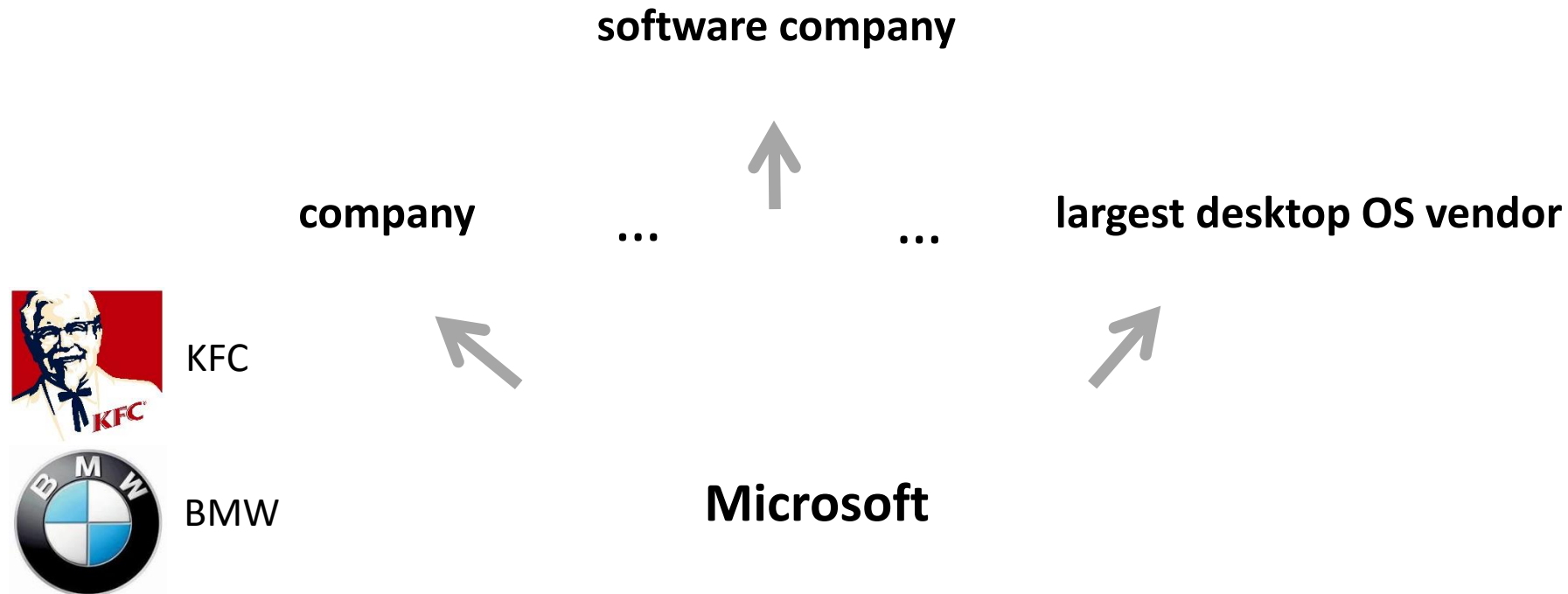




# Basic-level Conceptualization (BLC)

Category Level	Informative?	Distinctive?
Superordinate	No	Yes
Basic-level	Yes	Yes
Subordinate	Yes	No

Basic-level  
conceptualization



# Using $Rep(e, c)$ for BLC

- Our measure  $Rep(e, c) = P(c|e) * P(e|c)$  means:

Given  $e$ , the  $c$  should be its typical concept (**shortest distance**)

Given  $c$ , the  $e$  should be its typical entity (**shortest distance**)

A process of finding **concept nodes** having **shortest expected distance** with  $e$

- (With PMI) If we take the logarithm of our scoring function, we get:

$$\log Rep(e, c) = \log P(c|e) * P(e|c) = \log \frac{P(e, c)}{P(e)} * \frac{P(e, c)}{P(c)} = \log \frac{P(e, c)^2}{P(e)P(c)} = PMI(e, c) + \log P(e, c) = PMI^2$$

- (With Commute Time) The commute time between an **instance  $e$**  and a **concept  $c$**  is:

$$\begin{aligned} Time(e, c) &= \sum_{k=1}^{\infty} (2k) * P_k(e, c) = \sum_{k=1}^T (2k) * P_k(e, c) + \sum_{k=T+1}^{\infty} (2k) * P_k(e, c) \\ &\geq \sum_{k=1}^T (2k) * P_k(e, c) + 2(T+1) * (1 - \sum_{k=1}^T P_k(e, c)) = 4 - 2 * Rep(e, c) \end{aligned}$$

# Precision@K & NDCG@K

- Metrics

- $$Precision@K = \frac{\sum_{i=1}^K rel_i}{K}$$
 (for **correctness** of concepts)

- $$nDCG_K = \frac{rel_1 + \sum_{i=2}^K \frac{rel_i}{\log i}}{ideal\_rel_1 + \sum_{i=2}^K \frac{ideal\_rel_i}{\log i}}$$
 (for **ranking** of concepts)

- Results

Precision@K	1	2	3	5	10	15	20
MI(e)	0.769	0.692	0.705	0.685	0.719	0.705	0.690
<b>PMI<sup>3</sup>(e)</b>	<b>0.885</b>	0.769	0.756	0.800	0.754	<b>0.733</b>	<b>0.721</b>
NPMI(e)	0.692	0.692	0.667	0.638	0.627	0.610	0.610
Typicality P(c e)	0.462	0.577	0.603	0.577	0.569	0.564	0.556
Typicality P(e c)	0.500	0.462	0.526	0.523	0.523	0.510	0.521
<b>Rep(e)</b>	0.846	<b>0.865</b>	<b>0.872</b>	<b>0.862</b>	<b>0.758</b>	0.731	0.719

NDCG@K	1	2	3	5	10	15	20
MI(e)	0.516	0.531	0.519	0.531	0.562	0.574	0.594
PMI <sup>3</sup> (e)	0.725	0.664	0.652	0.660	0.628	0.631	0.646
NPMI(e)	0.599	0.597	0.579	0.554	0.540	0.539	0.549
Typicality P(c e)	0.297	0.380	0.409	0.422	0.438	0.446	0.461
Typicality P(e c)	0.401	0.386	0.396	0.398	0.401	0.410	0.428
<b>Rep(e)</b>	<b>0.758</b>	<b>0.771</b>	<b>0.745</b>	<b>0.723</b>	<b>0.656</b>	<b>0.647</b>	<b>0.661</b>

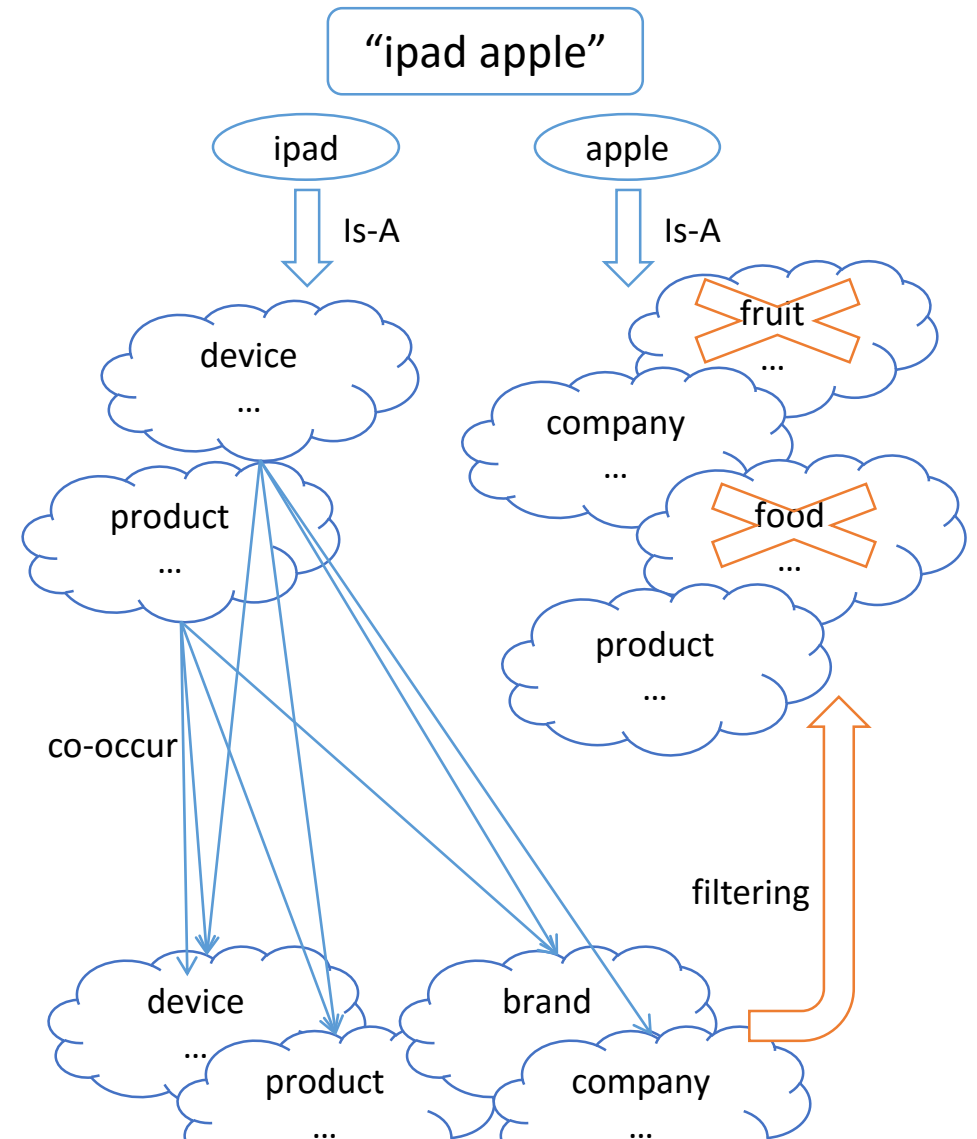
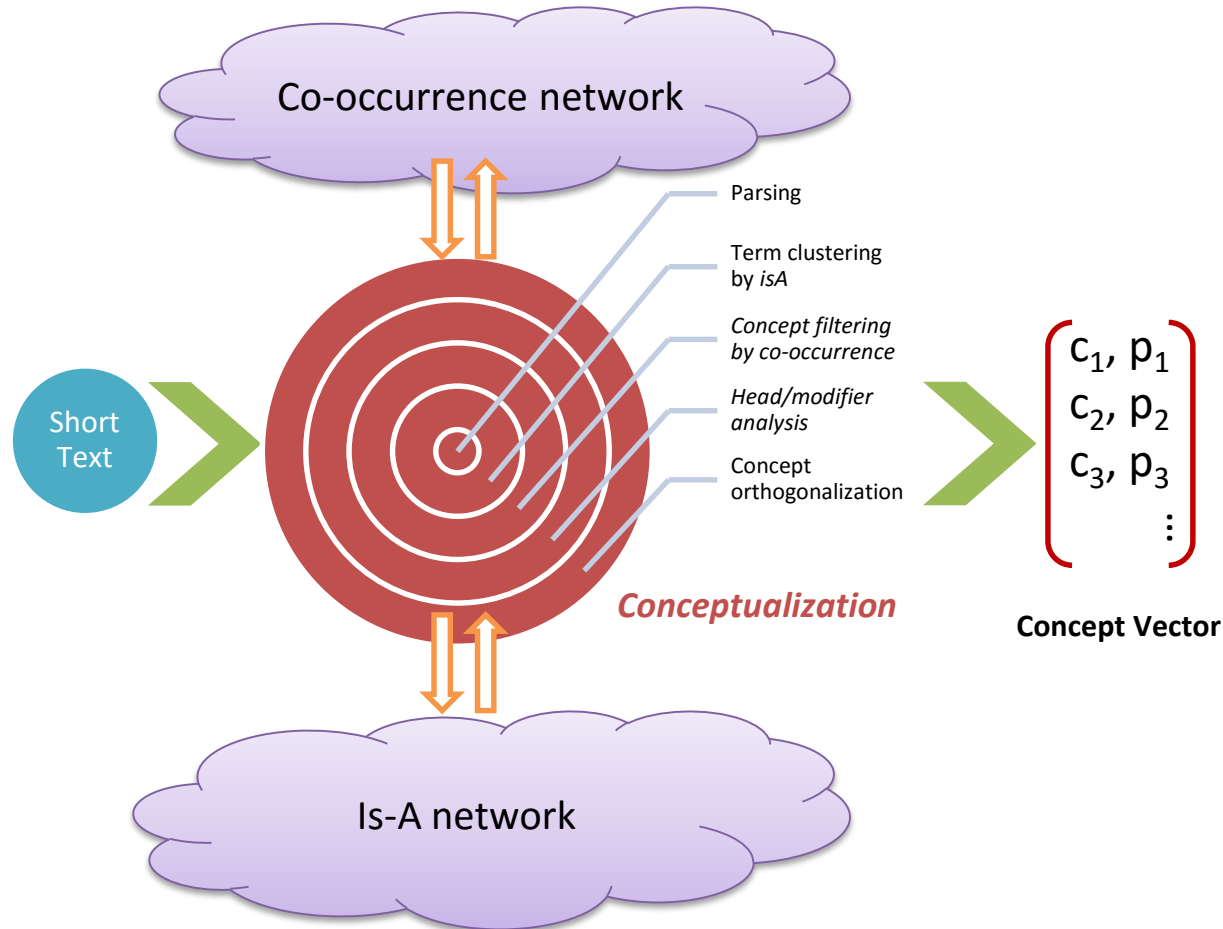
- Overall, our measure *Rep* performs well in both *Precision* and *NDCG*.
- Most important, it's well interpreted in theory

If the short text has context for the instance...

“Python Tutorial”

- Wen Hua, Zhongyuan Wang, Haixun Wang, Kai Zheng, and Xiaofang Zhou, [Short Text Understanding Through Lexical-Semantic Analysis](#), in *International Conference on Data Engineering (ICDE)*, April 2015. (**Best Paper Award**)
- Dongwoo Kim, Haixun Wang, and Alice Oh, [Context-Dependent Conceptualization](#), in *IJCAI*, 2013.
- Yangqiu Song, Haixun Wang, Zhongyuan Wang, Hongsong Li, and Weizhu Chen, [Short Text Conceptualization using a Probabilistic Knowledgebase](#), in *IJCAI*, 2011.

# Conceptualization Framework (ICDE 2015 Best Paper)

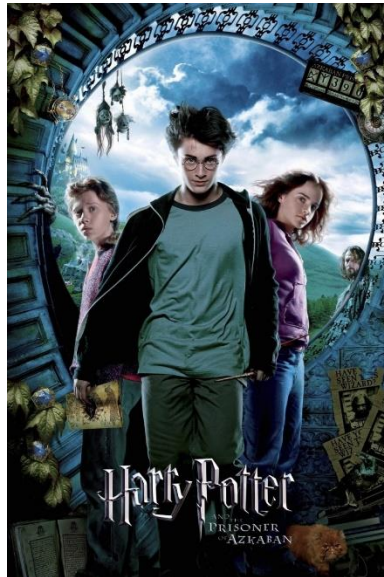


**If the short text contains verb, adjective...**

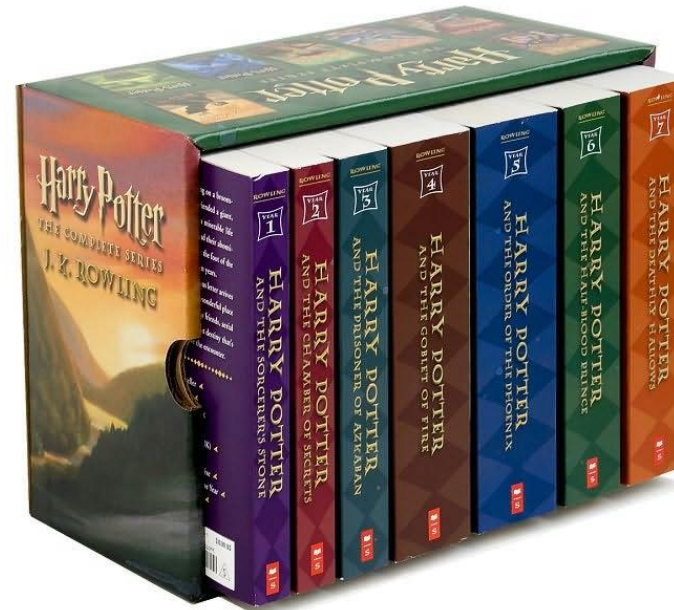
**“Dangerous Python”**

- Zhongyuan Wang, Kejun Zhao, Haixun Wang, Xiaofeng Meng, and Ji-Rong Wen, Query Understanding through Knowledge-Based Conceptualization, in *IJCAI*, 2015.

- Watch Harry Potter
- Read Harry Potter



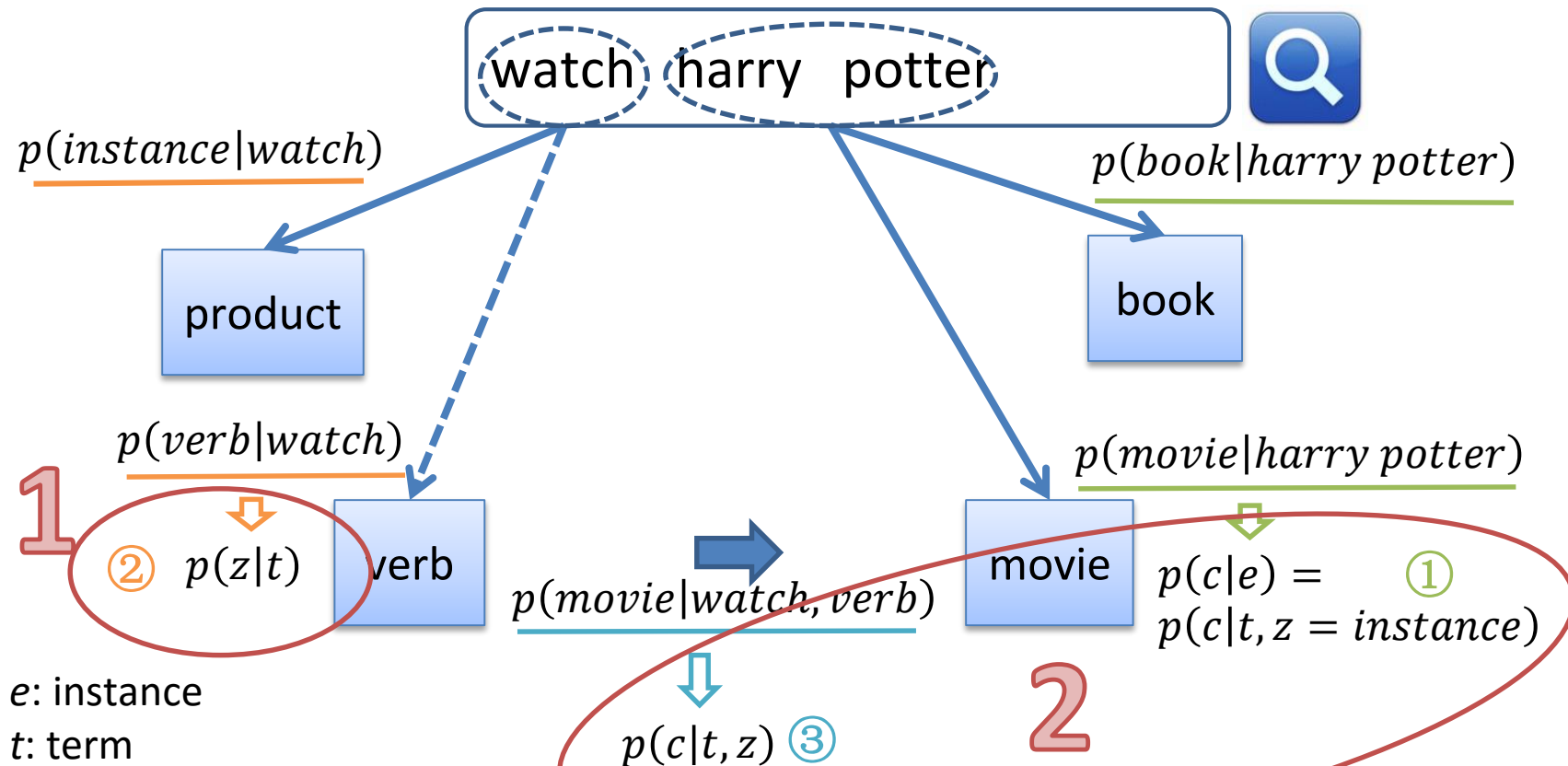
**Movie**



**Book**

# Mining Lexical Relationships

- Lexical knowledge represented by the probabilities





# Deriving Probabilities

- **Deriving  $p(\mathbf{z}|\mathbf{t})$ :**  $p(\mathbf{z}|\mathbf{t}) = \frac{n(\mathbf{t},\mathbf{z})}{n(\mathbf{t})}$

- **Deriving  $P(\mathbf{c}|\mathbf{t}, \mathbf{z})$**

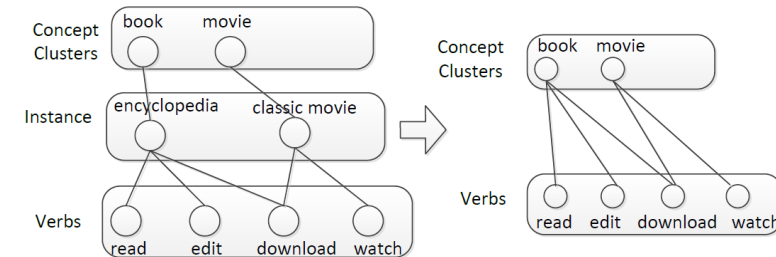
- Case 1:  $\mathbf{z}=\text{instance}$   $P(\mathbf{c}|\mathbf{t}, \mathbf{z} = \text{instance}) = p(\mathbf{c}|\mathbf{e})$
- Case 2:  $\mathbf{z}=\text{attribute}$   $P(\mathbf{c}|\mathbf{t}, \mathbf{z} = \text{attribute}) = p(\mathbf{c}|\mathbf{a})$

- Case 3:  $\mathbf{z}=\text{verb}$

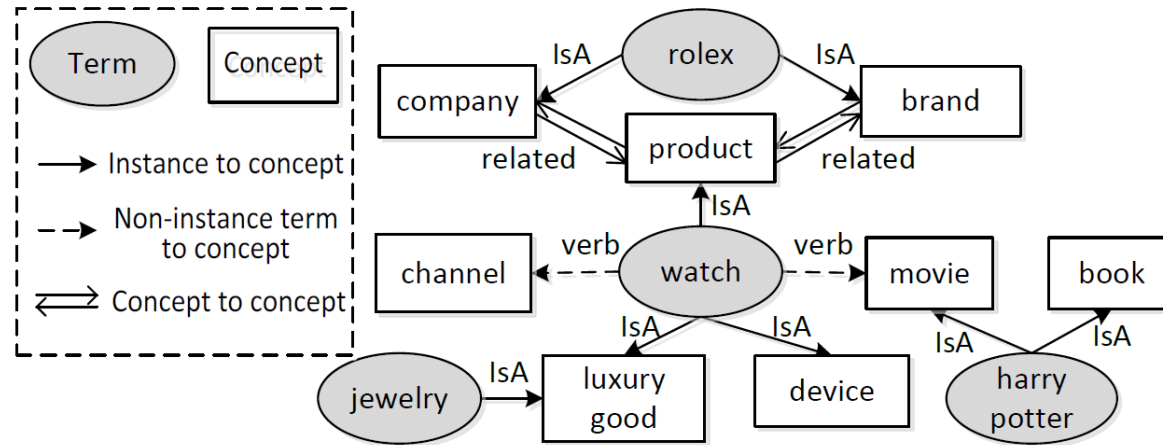
$$P(\mathbf{c}|\mathbf{t}, \mathbf{z} = \text{verb}) = \sum_{e \in \mathbf{c}} p(e, \mathbf{c}|\mathbf{t}, \mathbf{z} = \text{verb}) = \sum_{e \in \mathbf{c}} p(\mathbf{c}|\mathbf{e}) \times p(\mathbf{e}|\mathbf{t}, \mathbf{z} = \text{verb})$$

- Case 4:  $\mathbf{z}=\text{adjective}$

$$P(\mathbf{c}|\mathbf{t}, \mathbf{z} = \text{adjective}) = \sum_{e \in \mathbf{c}} p(\mathbf{c}|\mathbf{e}) \times p(\mathbf{e}|\mathbf{t}, \mathbf{z} = \text{adjective})$$



# Constructing an offline semantic network

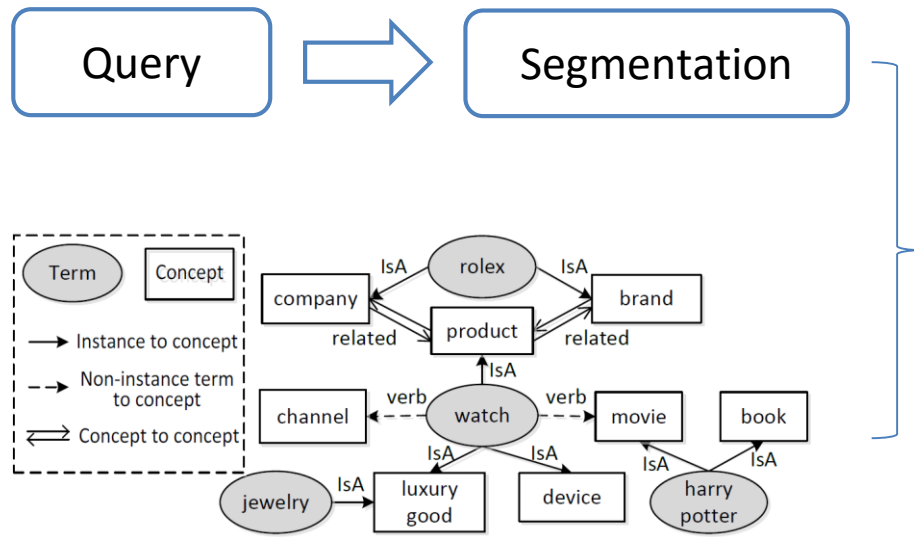


$$P(c|t) = \sum_z p(c|t, z) \times p(z|t)$$

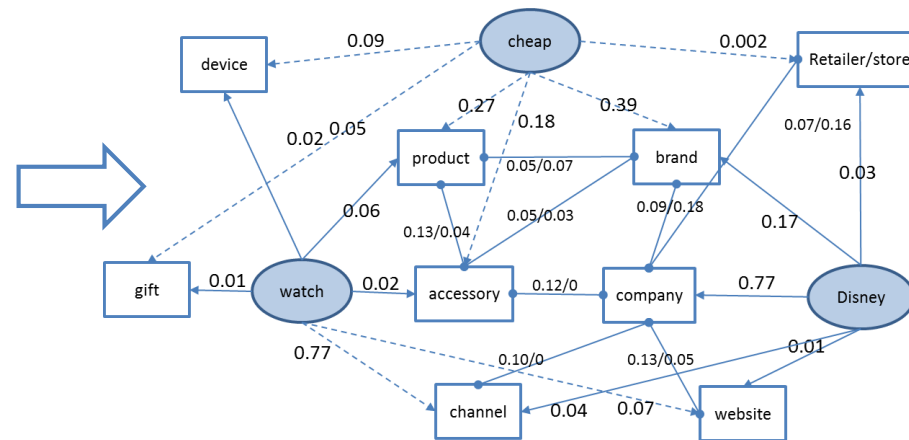
$$P(c_2|c_1) = \frac{\sum_{e_i \in c_1, e_j \in c_2} n(e_i, e_j)}{\sum_c \sum_{e_i \in c_1, e_j \in c} n(e_i, e_j)}$$

# Understanding Queries

- **Goal:** to rank the concepts and find:  
$$\arg \max_c p(c|t, q)$$



The offline semantic network



Random walk with restart [Sun et al., 2005]  
on the online subgraph

If the short text contains multiple instance...

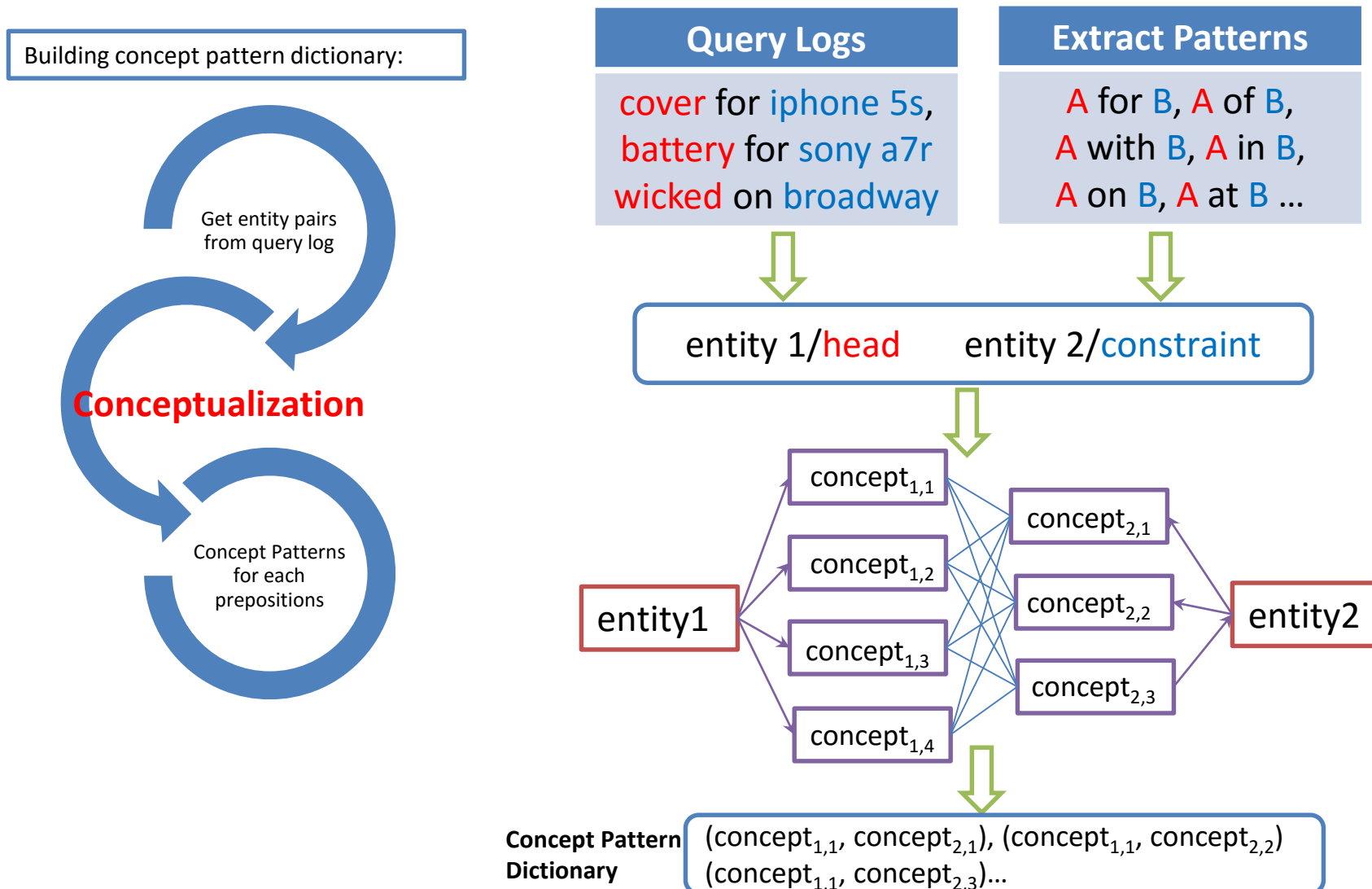
“DNN Tool Python”

Head

Modifier

- Zhongyuan Wang, Haixun Wang, and Zhirui Hu, [Head, Modifier, and Constraint Detection in Short Texts](#), in *International Conference on Data Engineering (ICDE)*, 2014.

# Mining Concept Patterns



# Why Concepts Can't Be Too General

- It may **cause too many concept pattern conflicts**: can't distinguish head and modifier for general concept pairs

	Head	Modifier
Derived Concept Pattern	device	company
Supporting Entity Pairs	iphone 4	verizon
	modem	comcast
	wireless router	comcast
	iphone 4	tmobile

	Head	Modifier
Derived Concept Pattern	company	device
Supporting Entity Pairs	amazon books	kindle
	netflix	touchpad
	skype	windows phone
	netflix	ps3

Conflict

# Why Concepts Can't Be Too Specific

- It may generate concepts **with less representation**

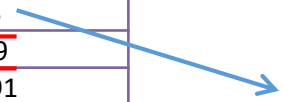
...	...
device	largest desktop OS vendor
device	largest software development company
device	largest global corporation
device	latest windows and office provider
...	...

- Concept level may **regress to entity level**
  - Large storage space: up to (million \* million) patterns

**We should use Basic-level  
Conceptualization (BLC)**

# Top Concept Patterns

Cluster size	Sum of Cluster Score	head;modifier;score
615	21146.91	breed;state;3572.98460224501
296	7752.357	game;platform;627.403476771856
153	3466.804	accessory;vehicle;533.93705094809
70	1182.59	browser;platform;132.612807637391
22	1010.993	requirement;school;271.407526294823
34	948.9159	drug;disease;154.602405333541
42	899.2995	cosmetic;skin condition;81.4659415003929
16	742.1599	job;city;279.03732555528
32	710.403	accessory;phone;246.513830851194
18	669.2376	software;platform;210.126322725878
20	644.4603	test;disease;239.774028397537
27	599.4205	clothes;breed;98.773996282851
19	591.3545	penalty;crime;200.544192793488
25	584.8804	tax;state;240.081818612579
16	546.5424	sauce;meat;183.592863621553
18	480.9389	credit card;country;142.919087972152
14	473.0792	food;holiday;145.54140330924
11	453.6199	mod;game;257.163856882439
29	435.0954	garment;sport;47.1533326845442
23	399.4886	career information;professional;73.2726483731257
15	386.065	song;instrument;128.189481818135
18	378.213	bait;fish;78.0426514113169
22	372.2948	study guide;book;50.8339765053921
19	340.8953	plugins;browser;55.0326072627126
14	330.5753	recipe;meat;88.2779863422951
18	321.4226	currency;country;110.825444188352
13	318.0272	lens;camera;186.081673263957
9	316.973	decoration;holiday;130.055844126533
16	314.875	food;animal;73.38544366514

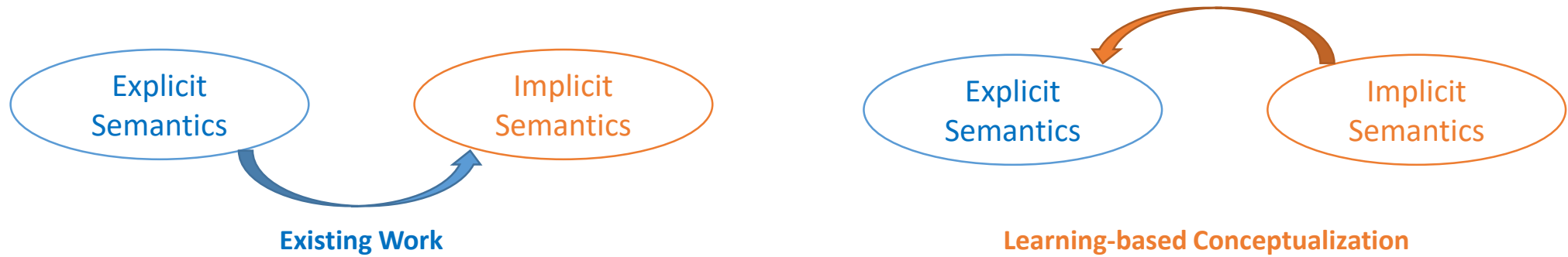


game	platform
game	device
video game	platform
game	console game pad
game	gaming platform

↓ Detection

Game (Head)	Platform (Modifier)
angry birds	android
angry birds	ios
angry birds	windows 10
...	...

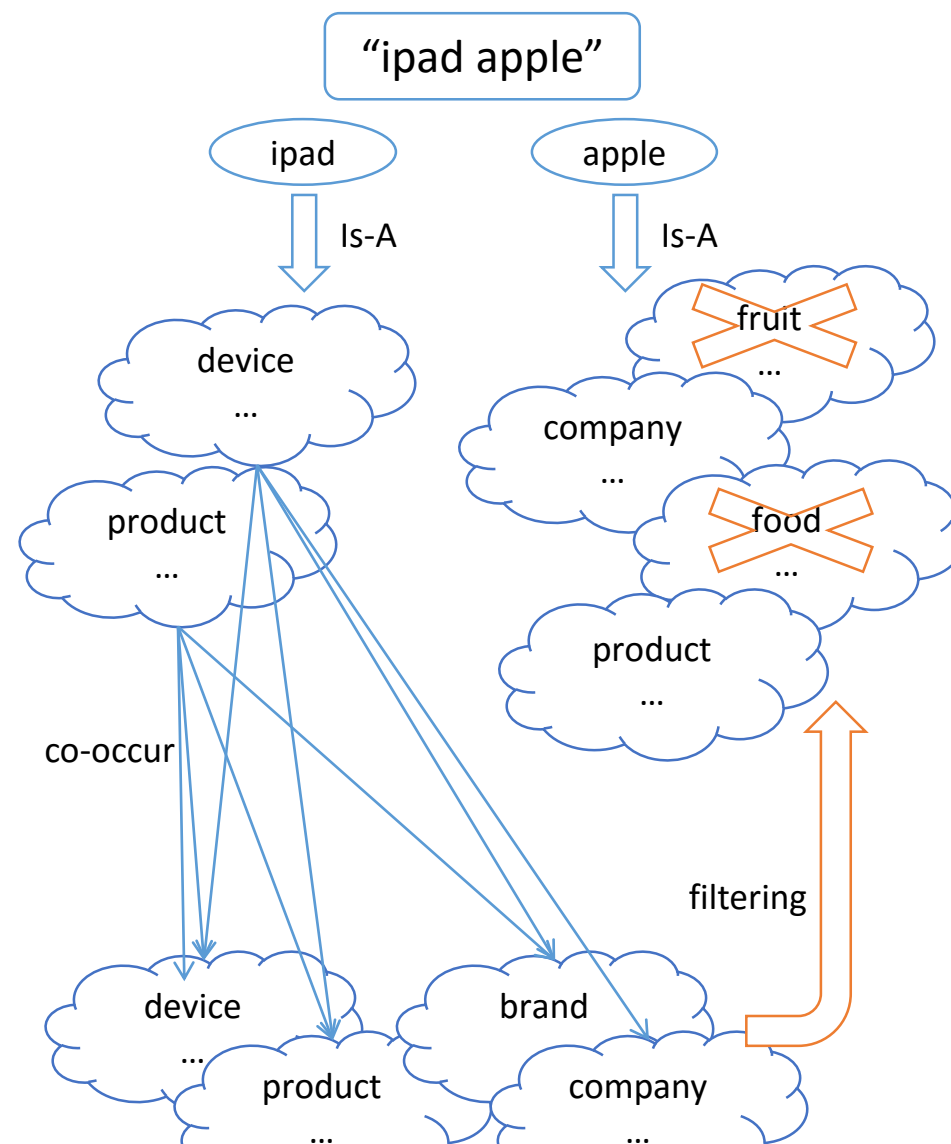
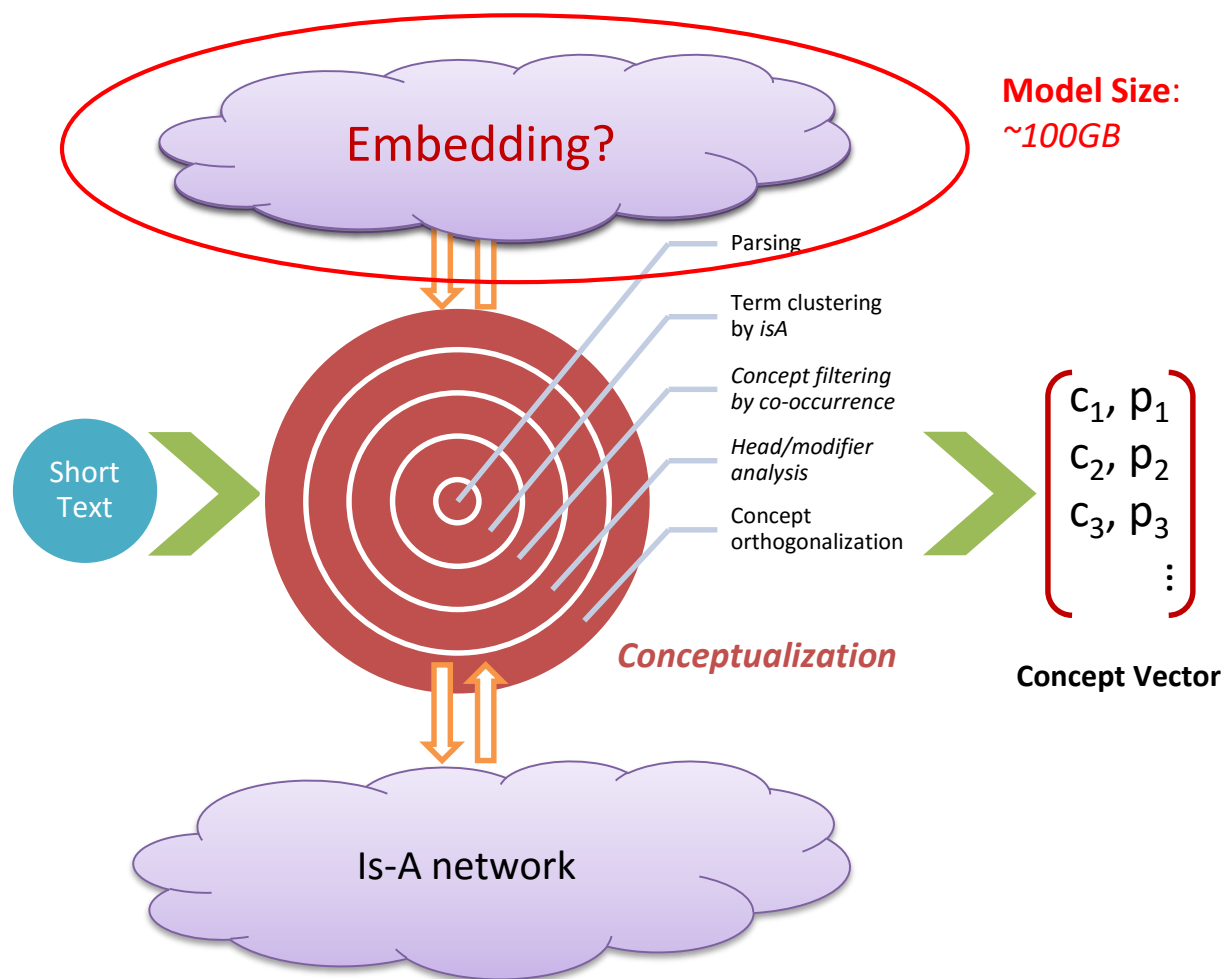




# Combine Explicit and Implicit Semantics (First Attempt): Learning-based Conceptualization

- Contextual Text Understanding in Distributional Semantic Space (CIKM2015)

# Previous Conceptualization Framework (ICDE 2015 Best Paper)



# Basic Idea of Learning-based Conceptualization

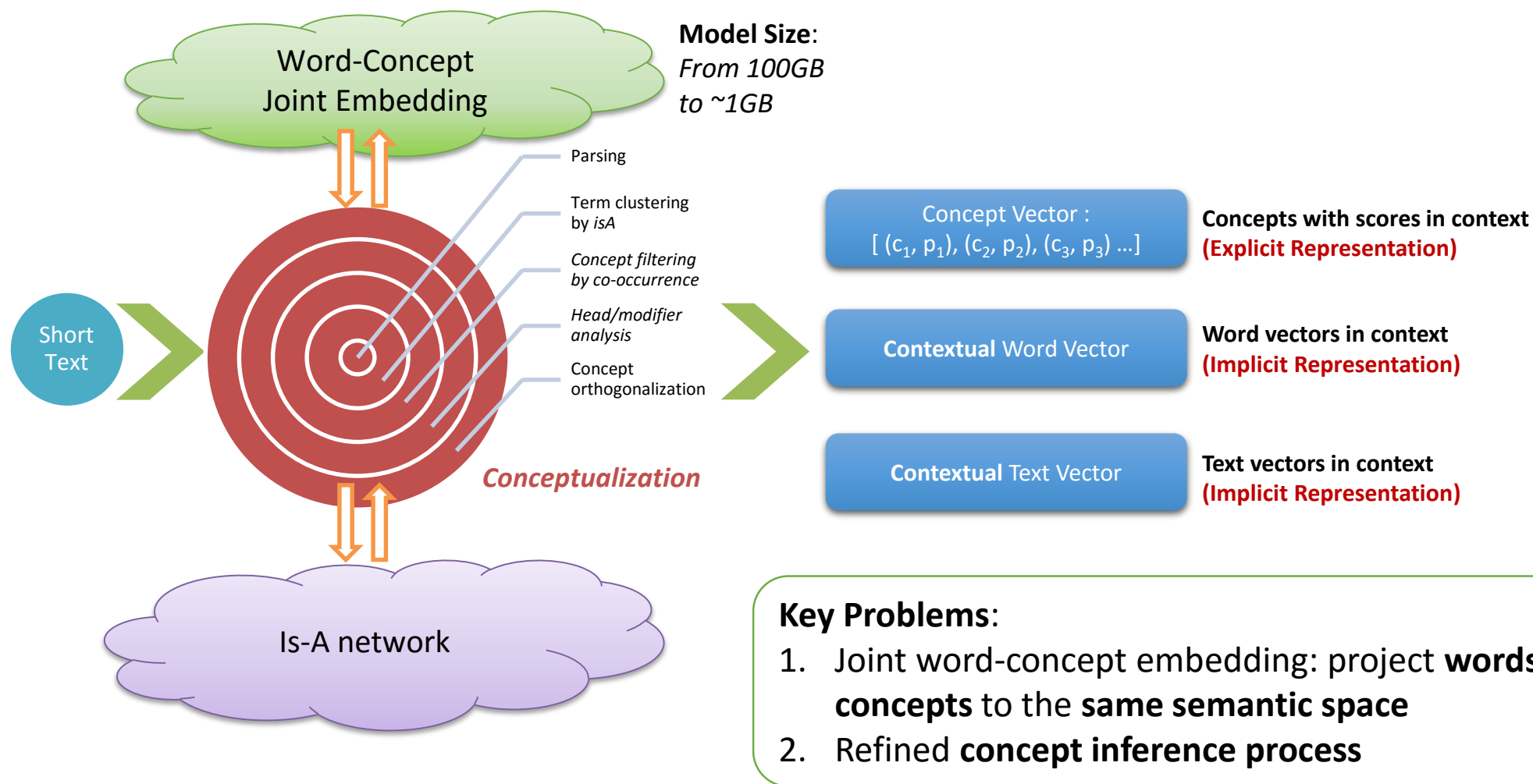
- **Target:** best representation for the short text:



- **Learning-based Conceptualization:** given a term  $t$ , with its context  $q$ , we want to find the probability of concept  $c$

$$f_{conceptualization}(c, t, q) = \underbrace{p(c|t)}_{\text{Traditional conceptualization based on Is-A network}} * \underbrace{p(c|q)}_{\text{Contextual knowledge based on Embedding}} = p(c|t) * \frac{\cos(\vec{q}, \vec{c})}{\sum_{c_i \in c(t)} \cos(\vec{q}, \vec{c}_i)}$$

# Learning-based Conceptualization Framework



# Joint Word-Concept Embedding

- Word embedding – Skip gram model

$$\sum_{i=1}^N \sum_{o=i-c}^{i+c} \log P(w_o | w_i)$$

- How to incorporate concept

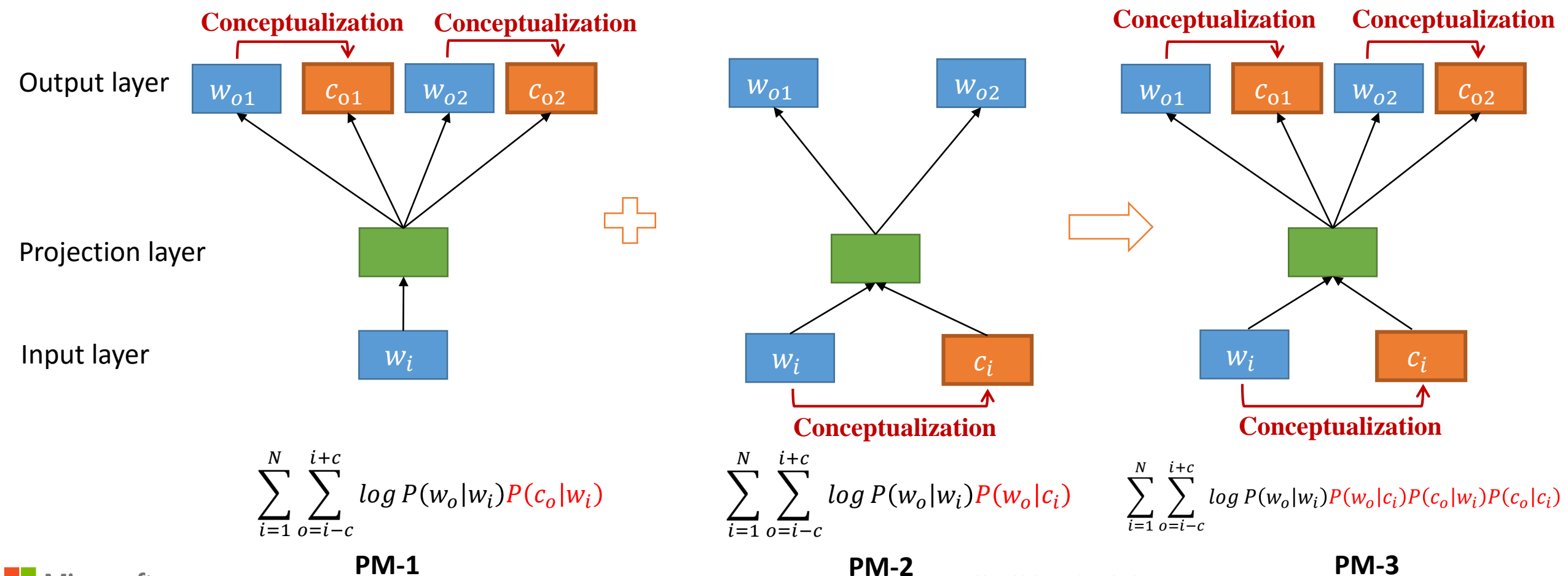
- Concept as input  $\sum_{i=1}^N \sum_{o=i-c}^{i+c} \log P(w_o | c_i)$

- Concept as output  $\sum_{i=1}^N \sum_{o=i-c}^{i+c} \log P(c_o | w_i)$

- Concept as both input and output  $\sum_{i=1}^N \sum_{o=i-c}^{i+c} \log P(c_o | c_i)$

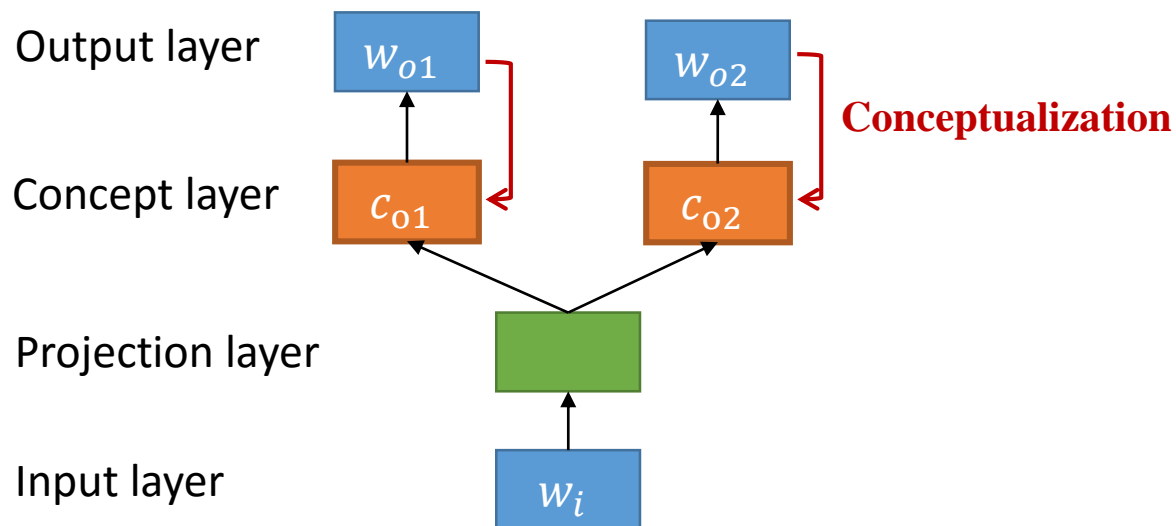
# Approach 1: Parallel Joint-Embedding Models

- Assume conditionally independent between the word and concept



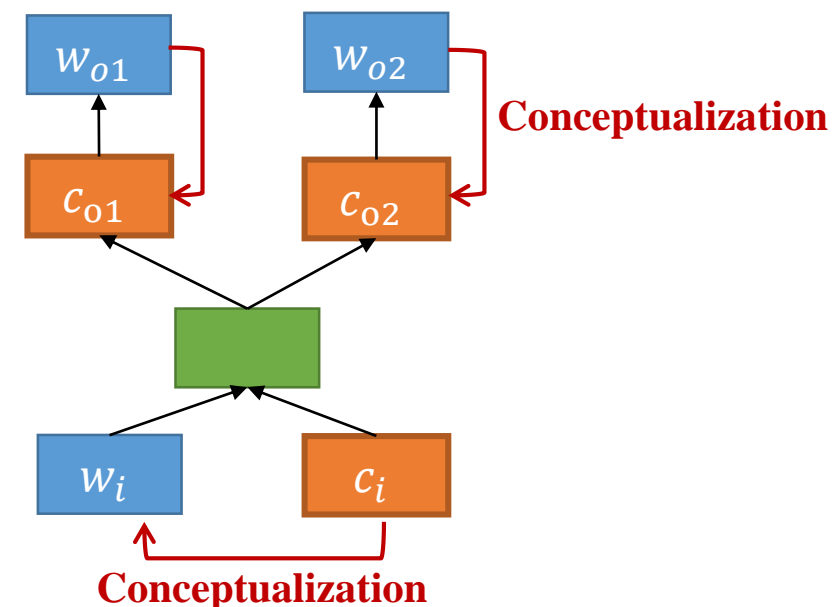
# Approach 2: Generative Joint-Embedding Models

- Assume conditionally dependent between output word and output concept: A word is selected by firstly select the class it belongs to.



$$\sum_{i=1}^N \sum_{t=i-c}^{i+c} \log P(w_o | w_i) = \sum_{i=1}^N \sum_{t=i-c}^{i+c} \log P(c_o | w_i) P(w_o | c_o)$$

GM-1



$$\sum_{i=1}^N \sum_{t=i-c}^{i+c} \log P(w_o | w_i, c_i) = \sum_{i=1}^N \sum_{t=i-c}^{i+c} \log P(c_o | w_i) P(c_o | c_i) P(w_o | c_o)$$

GM-2

# Evaluation: Word Similarity in Context

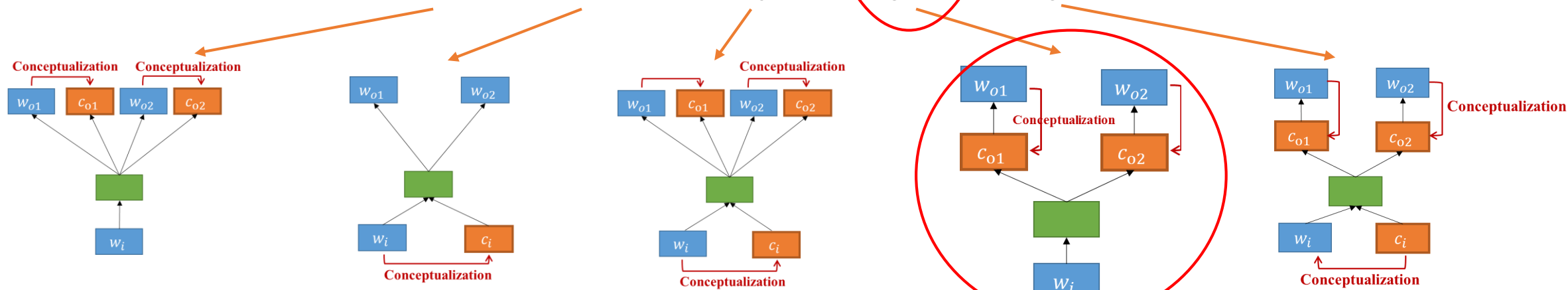
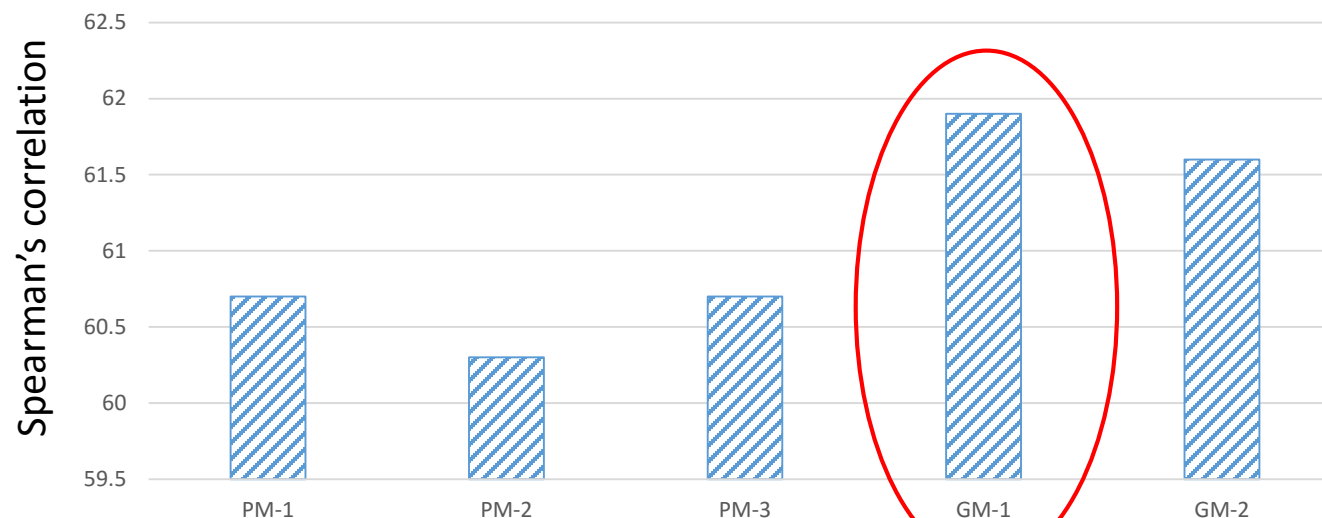
- Evaluation Setting
  - **Public dataset:** Eric Huang et al. (2012) Improving Word Representations via Global Context and Multiple Word Prototypes.
  - **Task:** Given *word1*, *word2* and their contexts, compute *similarity(word1, word2)*
- Metric: the **Spearman's correlation** between human judgement and embedding similarity score.

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}$$



# Word-Concept Embedding Evaluation Results

Non-contextual representations



- **PM:** Parallel Joint-Embedding Model

- **GM:** Generative Joint-Embedding Model

# Conceptualization with Word-Concept Embedding

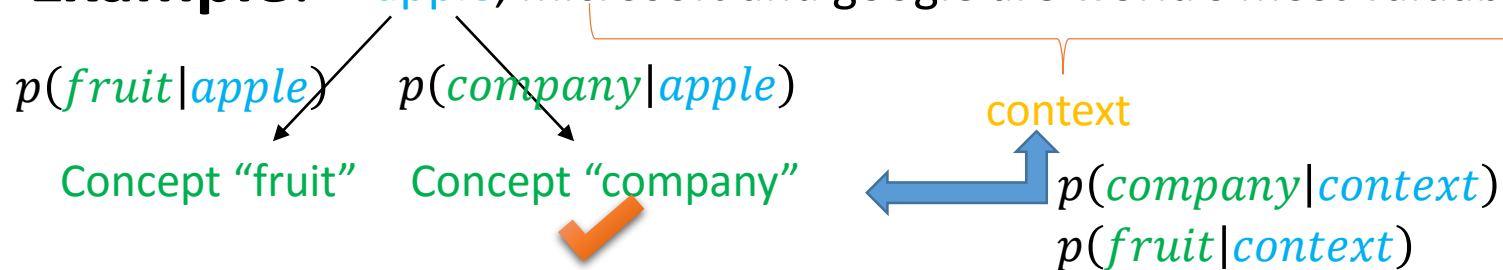
- **Concept Inference:** given a term  $t$ , with its context  $q$ , we want to find the probability of concept  $c$

$$f_{\text{conceptualization}}(c, t, q) = \underbrace{p(c|t)}_{\text{Traditional conceptualization based on Is-A network}} * \underbrace{p(c|q)}_{\text{Contextual knowledge based on Embedding}} = p(c|t) * \frac{\cos(\vec{q}, \vec{c})}{\sum_{c_i \in c(t)} \cos(\vec{q}, \vec{c}_i)}$$

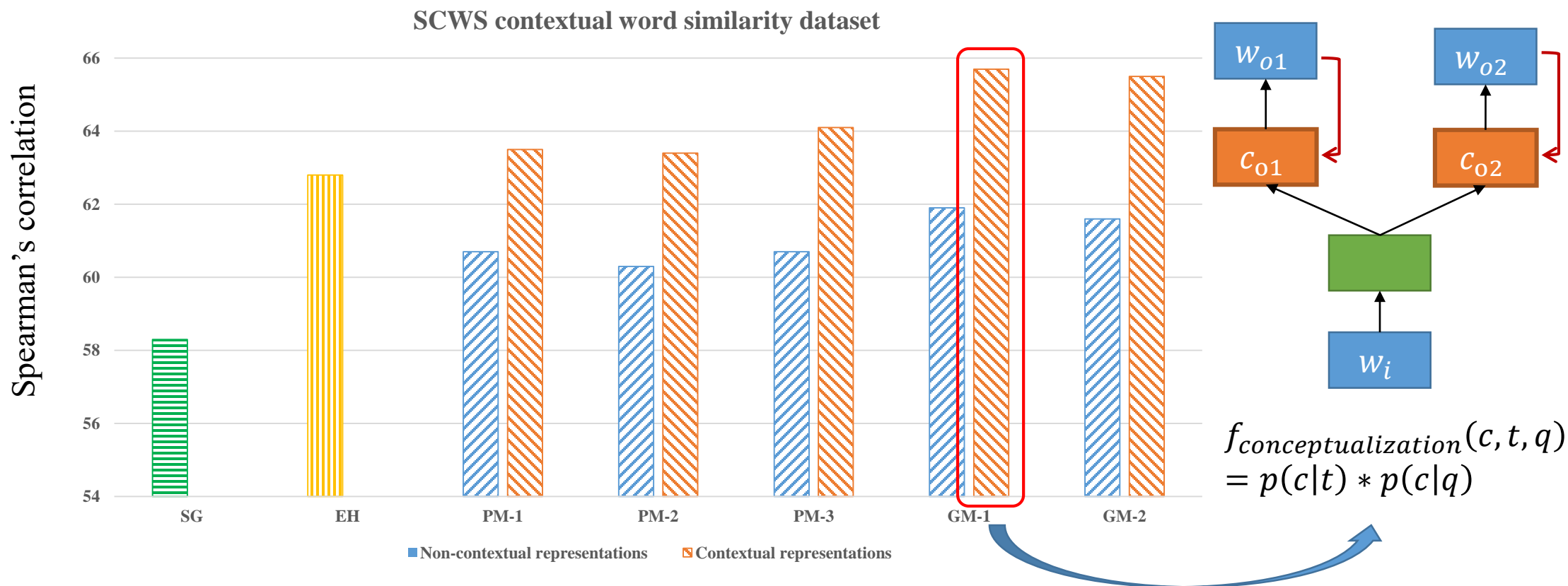
Traditional conceptualization  
based on Is-A network

Contextual knowledge  
based on Embedding

- **Conceptualization:**  $p(c|t, q) = \frac{f_{\text{conceptualization}}(c, t, q)}{\sum_{c_i \in c(t)} f_{\text{conceptualization}}(c_i, t, q)}$
- **Example:** “apple, microsoft and google are world’s most valuable brands.”



# Learning-based Conceptualization Evaluation Results



- **SG**: Skip-Gram (Word2Vec)
- **EH**: Eric Huang's Sense Embedding

- **PM**: Parallel Joint-Embedding Model
- **GM**: Generative Joint-Embedding Model

# Applications

- Short text understanding
- Short text similarity
- Ads/search semantic match
- Q/A system
- Query recommendation based on channels and articles
- Web table understanding
- ...

~Thank You~

<http://research.microsoft.com/probase/>

Contact: **Zhongyuan Wang**  
(email: zhy.wang # microsoft.com )