

# Context- and Content-aware Embeddings for Query Rewriting in Sponsored Search

Mihajlo Grbovic, Nemanja Djuric, Vladan Radosavljevic, [Fabrizio Silvestri](#), Narayan Bhamidipati

*Yahoo Labs - US & London*

# Query Rewriting

Try our state-of-the-art system for query rewriting.

bank of america loans	
Result	Relevance
bank_of_america_loans_personal	73.3375%
bank_of_america_personal_loans	70.0765%
bank_of_america_car_loans	68.5054%
bank_of_america_home_loans	68.2654%
bank_of_america_auto_loans	65.7220%
bank_of_america_mortgage	65.6517%
bank_of_america_personal_loans_online	64.9183%
bank_of_america_loan	64.1622%
bank_of_america_personal_loan	62.3031%
chase_personal_loans	62.23%
bank_of_america_student_loans	61.3762%

Input Type:

Output Type:

# Query Rewriting

Try our state-of-the-art system for query rewriting.

bank of america loans	
Result	Relevance
bank_of_america_loans_personal	73.3375%
bank_of_america_personal_loans	70.0765%
bank_of_america_car_loans	68.5054%
bank_of_america_home_loans	68.2654%
bank_of_america_auto_loans	65.7220%
bank_of_america_mortgage	65.6517%
bank_of_america_personal_loans_online	64.9183%
bank_of_america_loan	64.1622%
bank_of_america_personal_loan	62.3031%
chase_personal_loans	62.23%
bank_of_america_student_loans	61.3762%

Input Type:

Output Type:

# Query Rewriting for Advertising

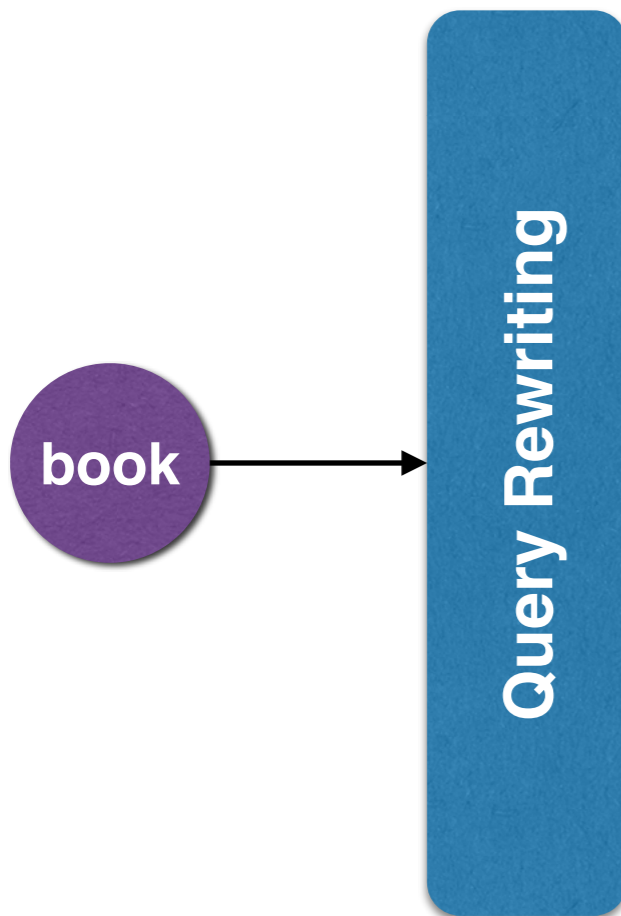


*bidterms* [Books at Amazon - Free 2-Day Shipping w/ Amazon Prime - Amazon.com](https://www.amazon.com/Books)  
[www.amazon.com/Books](https://www.amazon.com/Books)  
Low Prices on Millions of **Books**  
★★★★★ rating for amazon.com

*bidterms* [Cheap Audio Books](https://www.theworks.co.uk/audio-books)  
[TheWorks.co.uk/audio-books](https://www.theworks.co.uk/audio-books) Ad  
4.5 ★★★★★ user rating for theworks.co.uk  
Audio **Books** : Save at The Works Discount Store

*bidterms* [Kobo™ Free eBooks](https://www.kobobooks.com)  
[KoboBooks.com](https://www.kobobooks.com) Ad  
Millions Of Bestsellers And Free eBooks. Read On Almost Any Device!

# Query Rewriting for Advertising

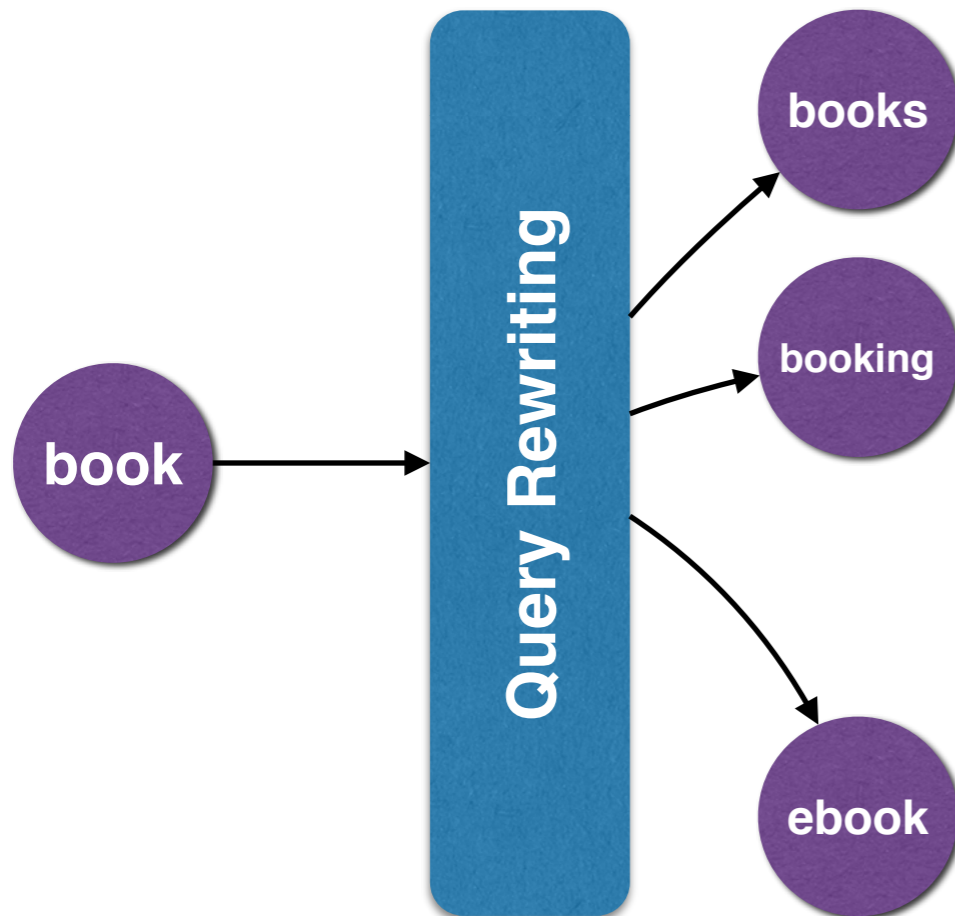


*bidterms* [Books at Amazon - Free 2-Day Shipping w/ Amazon Prime - Amazon.com](#)  
[www.amazon.com/Books](http://www.amazon.com/Books)  
Low Prices on Millions of **Books**  
★★★★★ rating for amazon.com

*bidterms* [Cheap Audio Books](#)  
[TheWorks.co.uk/audio-books](http://TheWorks.co.uk/audio-books) Ad  
4.5 ★★★★★ user rating for theworks.co.uk  
Audio **Books** : Save at The Works Discount Store

*bidterms* [Kobo™ Free eBooks](#)  
[KoboBooks.com](http://KoboBooks.com) Ad  
Millions Of Bestsellers And Free eBooks. Read On Almost Any Device!

# Query Rewriting for Advertising



*bidterms* [Books at Amazon - Free 2-Day Shipping w/ Amazon Prime - Amazon.com](#)  
[www.amazon.com/Books](http://www.amazon.com/Books)  
Low Prices on Millions of **Books**  
★★★★★ rating for amazon.com

*bidterms* [Cheap Audio Books](#)  
[TheWorks.co.uk/audio-books](http://TheWorks.co.uk/audio-books) Ad  
4.5 ★★★★★ user rating for theworks.co.uk  
Audio **Books** : Save at The Works Discount Store

*bidterms* [Kobo™ Free eBooks](#)  
[KoboBooks.com](http://KoboBooks.com) Ad  
Millions Of Bestsellers And Free eBooks. Read On Almost Any Device!

# (Lots of) Previous Work

- Graph-based Methods
  - Click Graph
  - Query-Flow Graph
  - Term-Query Graph
  - ...
- Clustering queries from history
- Supervised Learning-based methods
- ...

See more ads for:

[nike shoes](#)

[cheap nike shoes](#)

[nike shoes online](#)

[nike shoes sale](#)

[kids nike shoes](#)

Ads

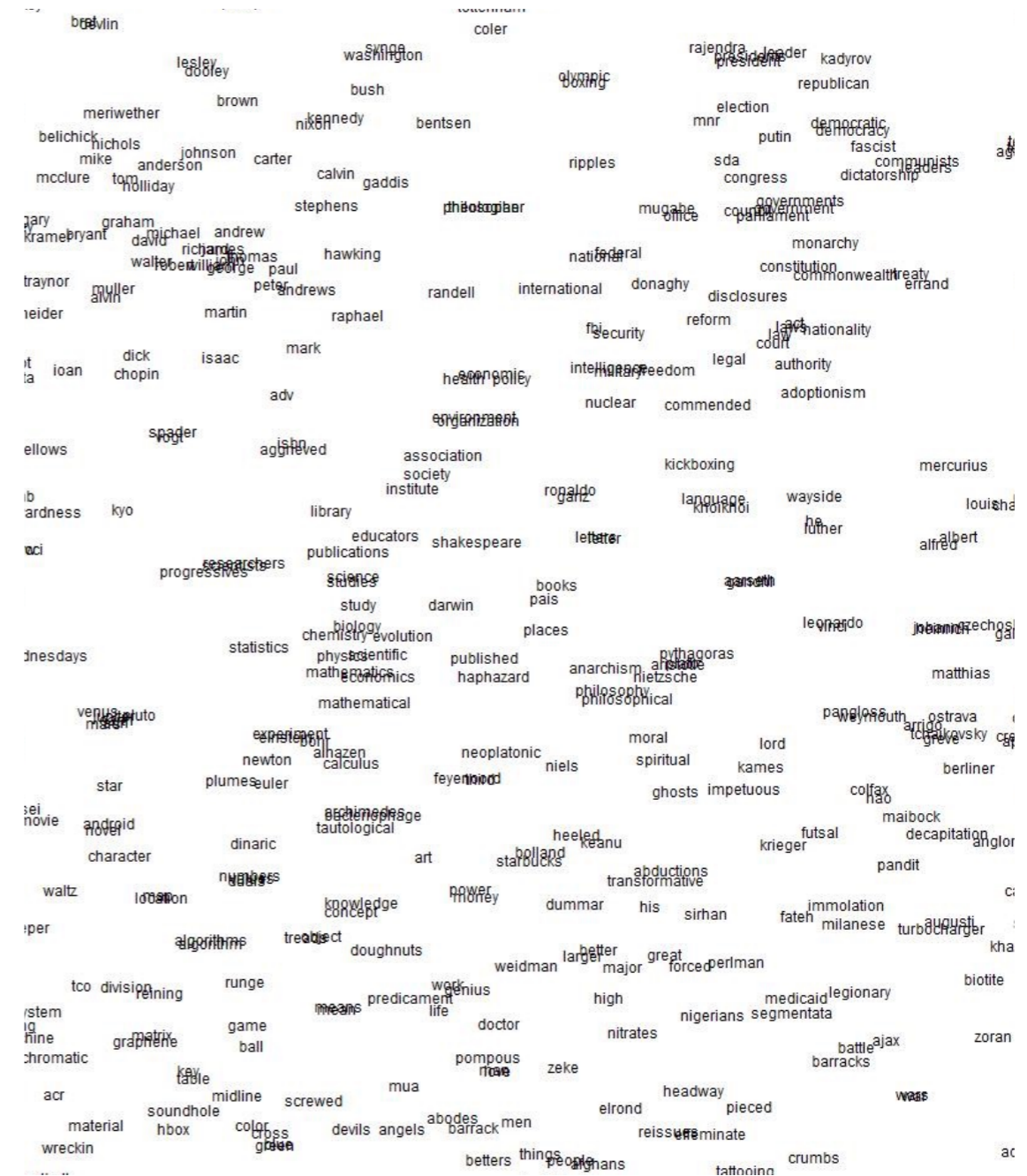
[\*\*Nike Shoes\*\*](#)

[lowpriceshopper.co](#)

**Best Nike Shoes Deals**

# Using Word Embeddings

- Sessions represent well the interactions of users with a search engines:
  - *queries*
  - *clicks*: algo results and ads
- We generate a distributed representation of actions in a common vector space:
  - inspired by *word2vec*





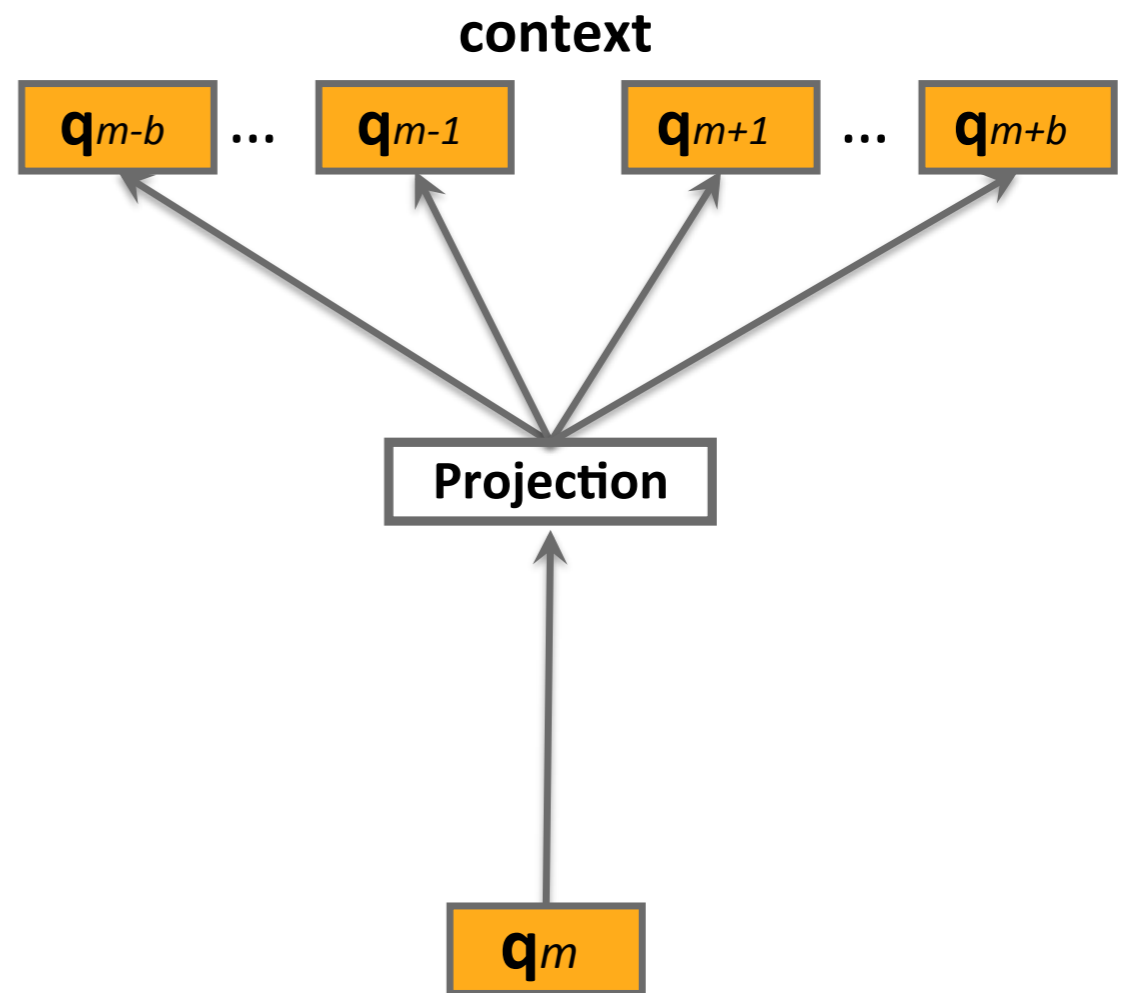
# Models for Embedding Interactions

- Context2Vec
- Content2Vec
- Context-Content2Vec



# Model 1: Context2vec

- Session as a sentence:
  - actions are words (e.g., 'nike shoes')
  - vectors are learned on each action/word

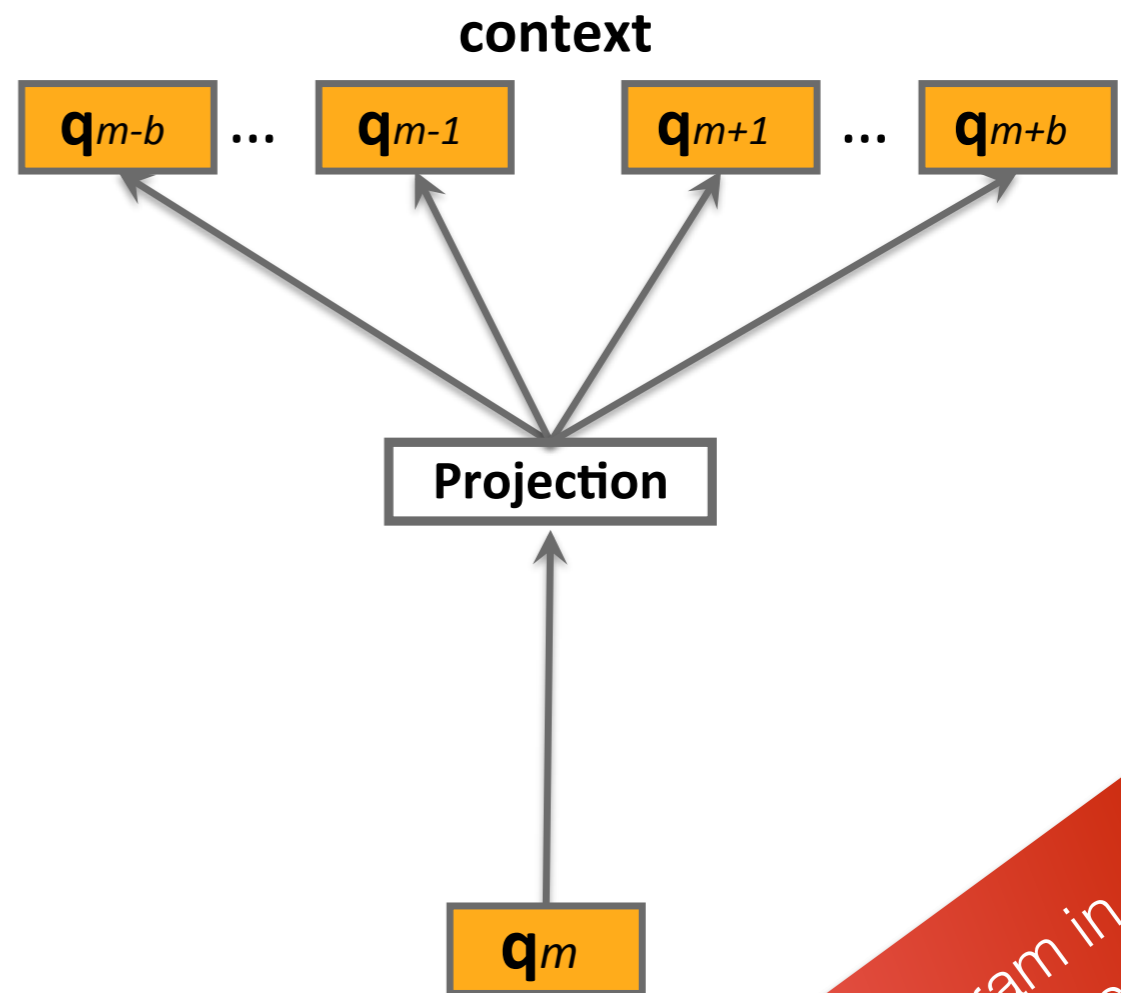


$$\mathcal{L} = \sum_{s \in \mathcal{S}} \sum_{q_m \in s} \sum_{-b \leq i \leq b, i \neq 0} \log \mathbb{P}(q_{m+i} | q_m).$$

$$\mathbb{P}(q_{m+i} | q_m) = \frac{\exp(\mathbf{v}_{q_m}^\top \mathbf{v}'_{q_{m+i}})}{\sum_{q=1}^V \exp(\mathbf{v}_{q_m}^\top \mathbf{v}'_q)}$$

# Model 1: Context2vec

- Session as a sentence:
  - actions are words (e.g., 'nike shoes')
  - vectors are learned on each action/word



$$\mathcal{L} = \sum_{s \in \mathcal{S}} \sum_{q_m \in s} \sum_{-b \leq i \leq b, i \neq 0} \log \mathbb{P}(q_{m+i} | q_m).$$

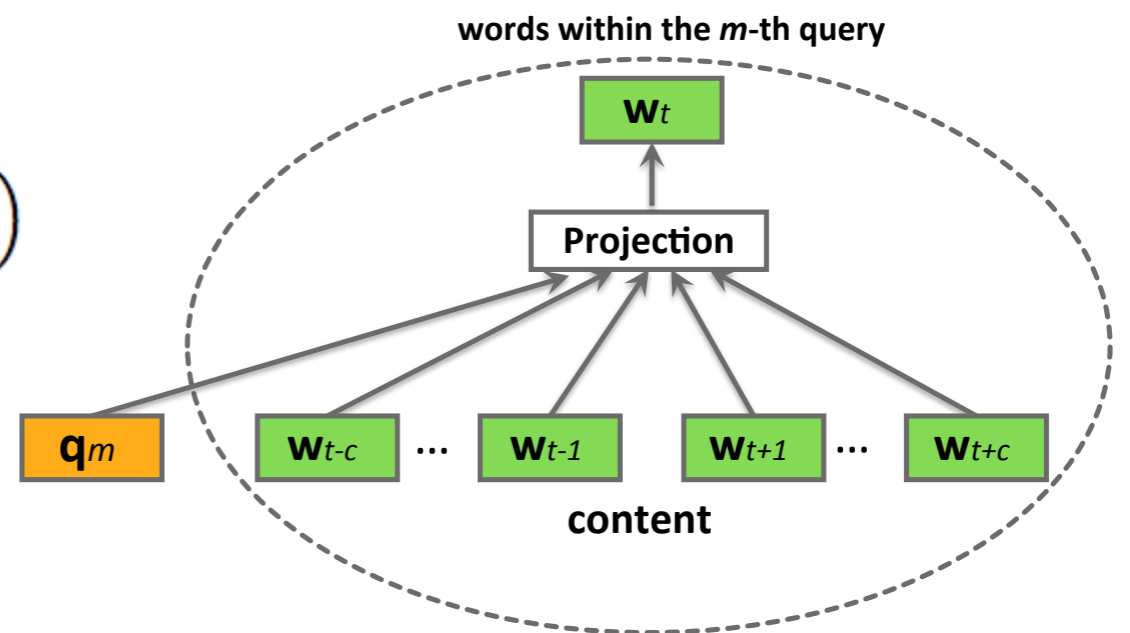
$$\mathbb{P}(q_{m+i} | q_m) = \frac{\exp(\mathbf{v}_{q_m}^\top \mathbf{v}'_{q_{m+i}})}{\sum_{q=1}^V \exp(\mathbf{v}_{q_m}^\top \mathbf{v}'_q)}$$

Skip-gram in word2vec

# Model 2: Content2Vec

- Consider the “content” of the action

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \left( \sum_{q_m \in s} \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) \right. \\ \left. + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$



# Model 2: Content2Vec

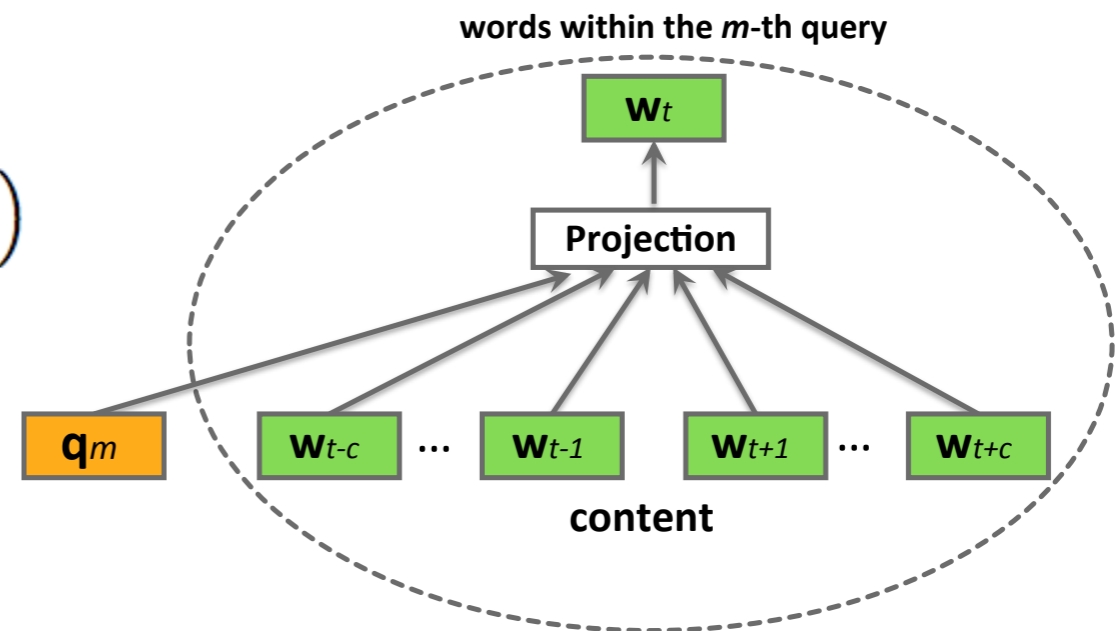
- Consider the “content” of the action

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \left( \sum_{q_m \in s} \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$

$$\mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) = \frac{\exp(\bar{\mathbf{v}}^\top \mathbf{v}'_{w_{mt}})}{\sum_{w=1}^V \exp(\bar{\mathbf{v}}^\top \mathbf{v}'_w)}$$

$$\bar{\mathbf{v}} = \frac{1}{2c+1} (\mathbf{v}_{q_m} + \sum_{-c \leq j \leq c, j \neq 0} \mathbf{v}_{w_{m,t+j}})$$

$$\mathbb{P}(q_m | w_{m1} : w_{mT_m}) = \frac{\exp(\bar{\mathbf{v}}_m^\top \mathbf{v}'_{q_m})}{\sum_{w=1}^V \exp(\bar{\mathbf{v}}_m^\top \mathbf{v}'_w)}$$



$$\bar{\mathbf{v}}_m = \frac{1}{T_m} \sum_{t=1}^{T_m} \mathbf{v}_{w_{mt}}$$

# Model 2: Content2Vec

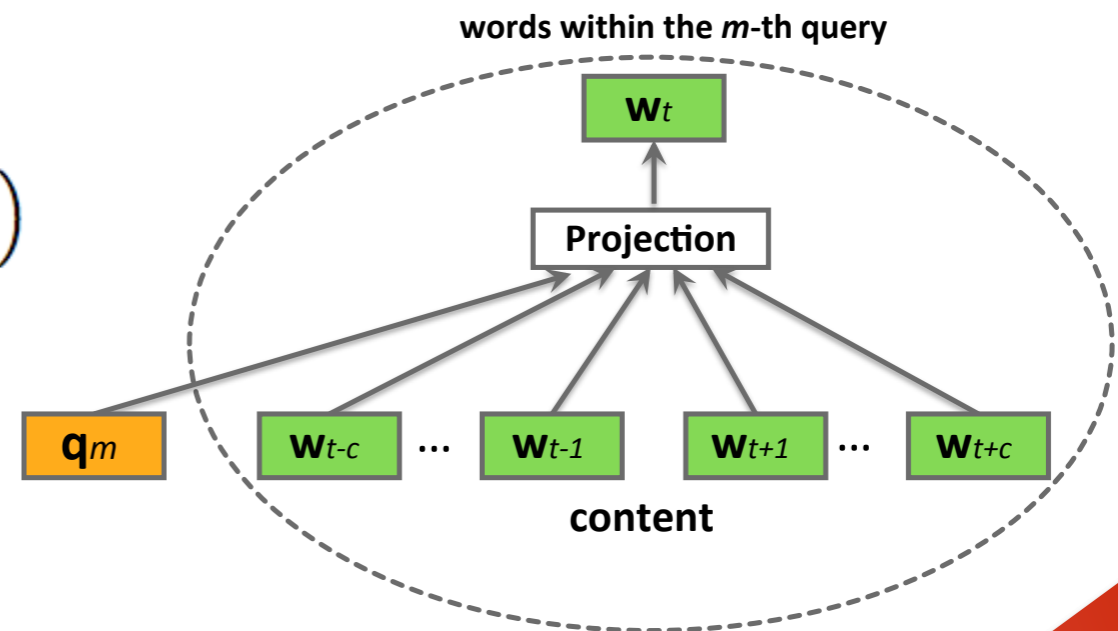
- Consider the “content” of the action

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \left( \sum_{q_m \in s} \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$

$$\mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) = \frac{\exp(\bar{\mathbf{v}}^\top \mathbf{v}'_{w_{mt}})}{\sum_{w=1}^V \exp(\bar{\mathbf{v}}^\top \mathbf{v}'_w)}$$

$$\bar{\mathbf{v}} = \frac{1}{2c+1} (\mathbf{v}_{q_m} + \sum_{-c \leq j \leq c, j \neq 0} \mathbf{v}_{w_{m,t+j}})$$

$$\mathbb{P}(q_m | w_{m1} : w_{mT_m}) = \frac{\exp(\bar{\mathbf{v}}_m^\top \mathbf{v}'_{q_m})}{\sum_{w=1}^V \exp(\bar{\mathbf{v}}_m^\top \mathbf{v}'_w)}$$



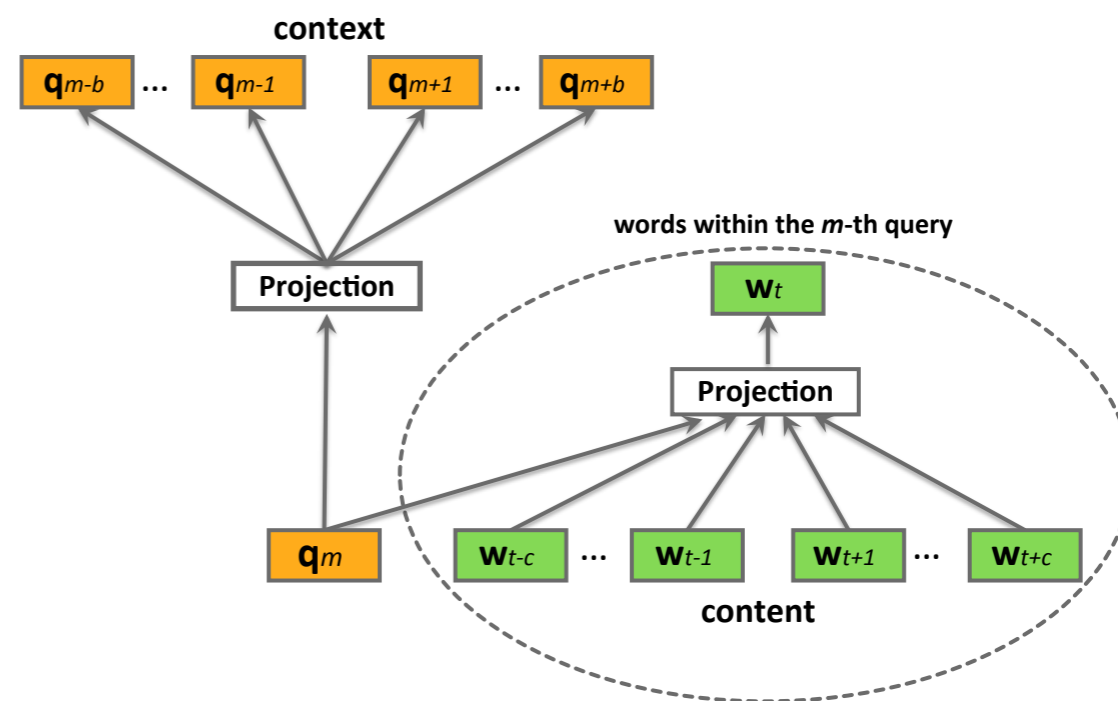
$$\bar{\mathbf{v}}_m = \frac{1}{T_m} \sum_{t=1}^{T_m} \mathbf{v}_{w_{mt}}$$

Similar to Paragraph2Vec

# Model 3: Context-Content2Vec

- A combination of the previous two models
- When an action is not popular we give more importance to its content

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \sum_{q_m \in s} \left( \sum_{-b \leq i \leq b, i \neq 0} \log \mathbb{P}(q_{m+i} | q_m) \right. \\ \left. + \alpha_m \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) \right. \\ \left. + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$



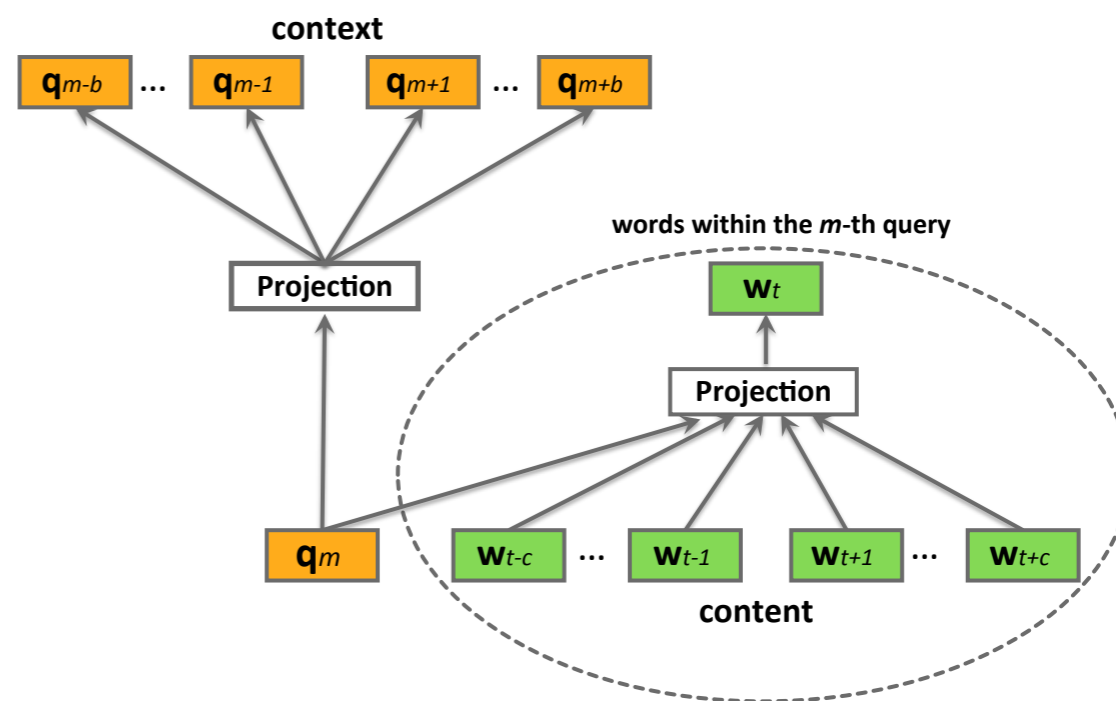
# Model 3: Context-Content2Vec

- A combination of the previous two models
- When an action is not popular we give more importance to its content

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \sum_{q_m \in s} \left( \sum_{-b \leq i \leq b, i \neq 0} \log \mathbb{P}(q_{m+i} | q_m) \right. \\ \left. + \alpha_m \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) \right. \\ \left. + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$

$$\alpha_m = \frac{1}{\log(1 + K_m)}$$

$K_m$  is the frequency of the  $K$ -th actions





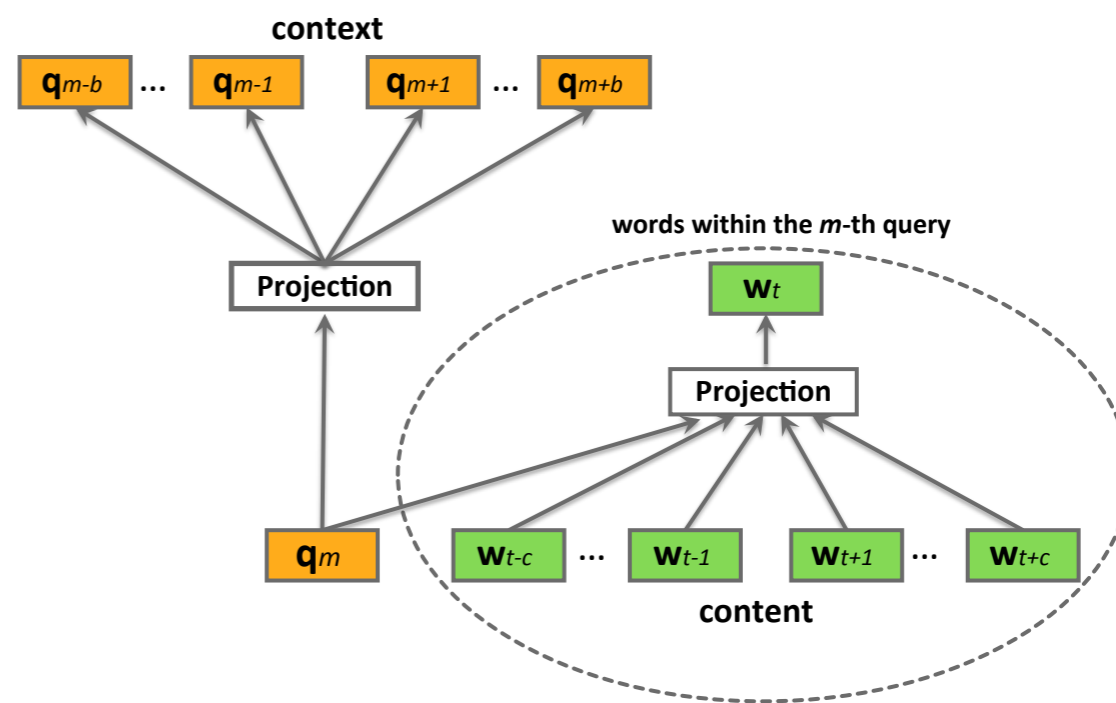
# Model 3: Context-Content2Vec

- A combination of the previous two models
- When an action is not popular we give more importance to its content

$$\mathcal{L} = \sum_{s \in \mathcal{S}} \sum_{q_m \in s} \left( \sum_{-b \leq i \leq b, i \neq 0} \log \mathbb{P}(q_{m+i} | q_m) \right. \\ \left. + \alpha_m \log \mathbb{P}(q_m | w_{m1} : w_{mT_m}) \right. \\ \left. + \sum_{w_{mt} \in q_m} \log \mathbb{P}(w_{mt} | w_{m,t-c} : w_{m,t+c}, q_m) \right)$$

$$\alpha_m = \frac{1}{\log(1 + K_m)}$$

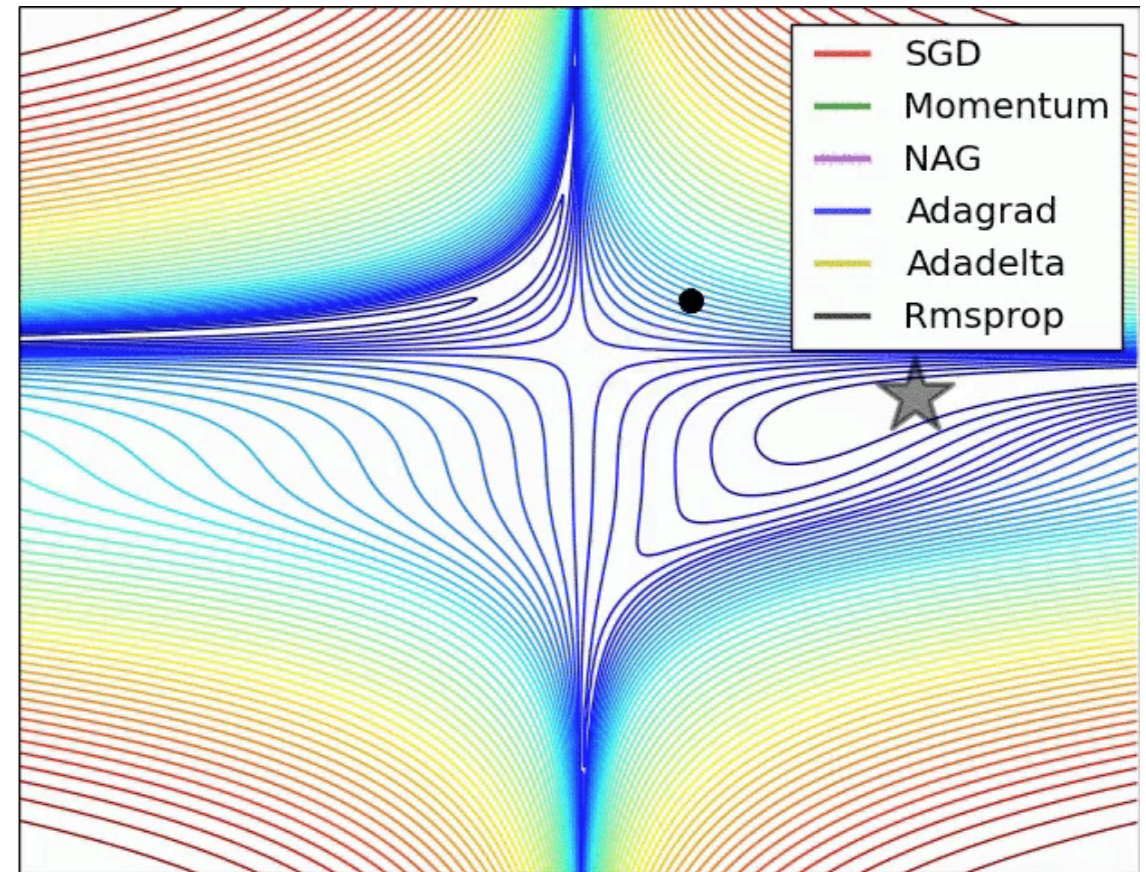
$K_m$  is the frequency of the K-th actions



Original Contribution

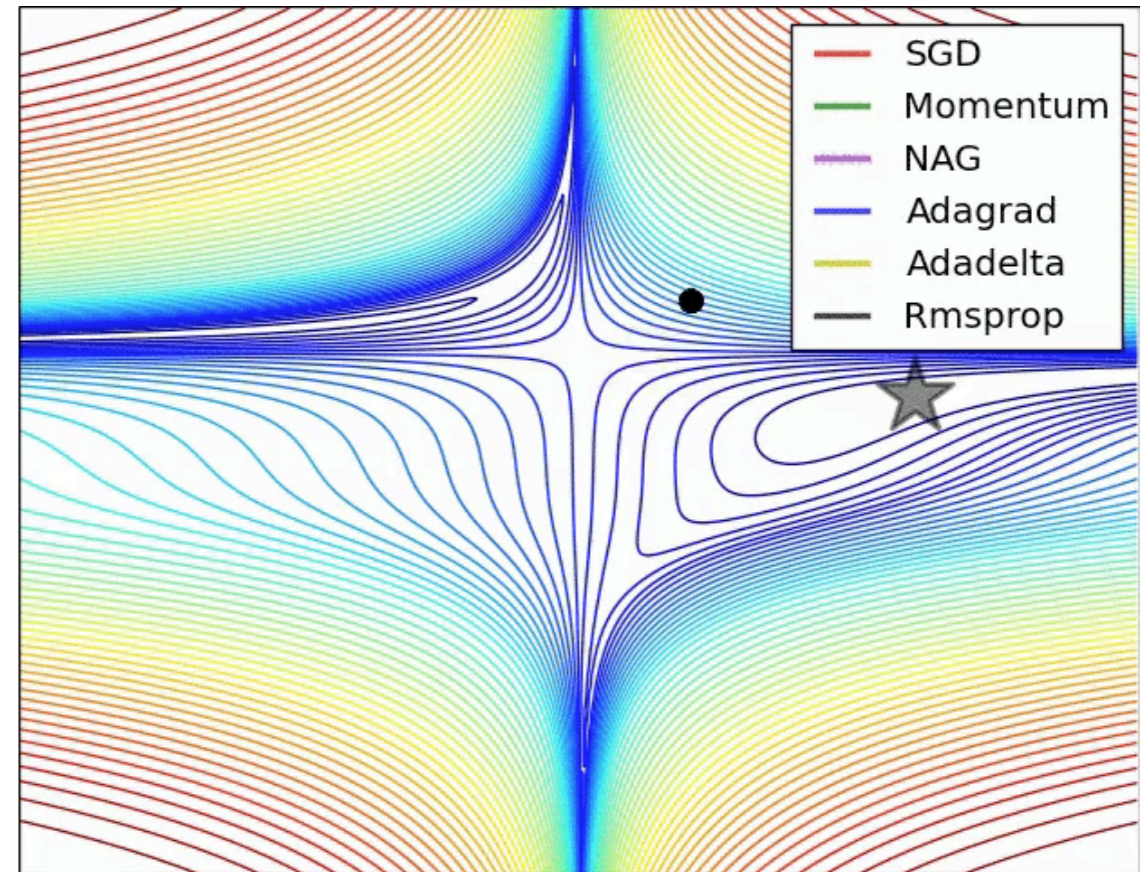
# Optimization

- Stochastic Gradient Ascent
- Negative Sampling:
  - 10 Negative samples
  - Negatives chosen in order to deal with clicks and skips
- Dimension  $D = 300$
- Context neighborhood = 5
- Content neighborhood = 7



# Optimization

- Stochastic Gradient Ascent
- Negative Sampling:
  - 10 Negative samples
  - Negatives chosen in order to deal with clicks and skips
- Dimension  $D = 300$
- Context neighborhood = 5
- Content neighborhood = 7



# Experiments

- Session Data
  - *12 billion* sessions collected on the US website of Yahoo Search
- Query Content Data
  - from the *45 million* most frequent queries
- Click Data
  - For all the sessions we collected the clicks on algo and sponsored results
- Query Flow Graph (QFG) as baseline



# Anecdotes

Query	cx2vec	cx-cn2vec
minnesota insurance exam crossword puzzles	satellite tv otego satellite tv menominee satellite tv west end satellite tv townsend satellite tv lake sara	minnesota insurance minnesota insurance license practice exams online insurance exam crossword puzzles colorado insurance exam crossword puzzles online minnesota insurance exam crossword puzzles
microwave food safety	staphylococcal enteritis definition salmonella enteritis definition listeria monocytogenes prevention e coli cdc preventing cross contamination	microwave oven food safety microwave baby food safety microwave food safety studies microwave food safety issues foodsafety.com
what to cook in cast iron skillet	steak sauce substitute montreal seasoning ingredients ground turkey breakfast sausage how to cook steak on cast iron skillet reseason cast iron	cast iron skillet recipes how to cook with cast iron skillet how to cook in a cast iron skillet how to cook with a cast iron skillet chicken in cast iron skillet
iphone 6 repair services	mp3attic music donar ovulos en elche credit a la consommation rapide smart phone repair service social security disability bronx ny	at&t iphone repair service iphone 5c repair service iphone repair iphone service repair iphone repair services

# The Effect of Ad Clicks

<b>makeup (no ads)</b>	<b>makeup (with ads)</b>
makeup tips	lipstick
fashion makeup	<b>mac makeup</b>
make up	<b>makeup sets</b>
makeup pictures	eye shadow
makeup images	<b>makeup covergirl</b>
makeup tutorial	makeup items

# The Effect of Ad Clicks

<b>makeup (no ads)</b>	<b>makeup (with ads)</b>	
makeup tips	lipstick	
fashion makeup	<b>mac makeup</b>	
make up	<b>makeup sets</b>	
makeup pictures	eye shadow	
<b>makeup</b>	<b>snowboarding (no ads)</b>	<b>snowboarding (with ads)</b>
makeup	snowbaording	snowboards
	snow boarding	<b>snowboarding gear</b>
	snowboarding information	<b>burton snowboarding</b>
	snowboarding jumps	<b>snowboard deals</b>
	snowboard pics	<b>snowboards on sale</b>
	shaun white snowboarding	snowboarding mountains

# The Effect of Ad Clicks

<b>makeup (no ads)</b>	<b>makeup (with ads)</b>
makeup tips	lipstick
fashion makeup	<b>mac makeup</b>
make up	<b>makeup sets</b>
makeup pictures	eye shadow
<b>snowboarding (no ads)</b>	<b>snowboarding (with ads)</b>
makeup snowboarding	snowboards
makeup snow boarding	<b>snowboarding gear</b>
snowboarding information	<b>burton snowboarding</b>
snowboarding jump	<b>seafood (no ads)</b>
snowboard pics	<b>seafood (with ads)</b>
shaun white snowbo	sea food
	<b>seafood restaurant</b>
	crab legs
	<b>seafood restaurants</b>
	best seafood
	crab shack
	<b>seafood market</b>
	oysters
	sea food
	<b>seafood market</b>
	<b>sea food menu</b>



# A Broader Comparison

Original	QFG <sub>ad+link</sub>	cn2vec	cx-cn2vec	cx-cn2vec <sub>ad</sub>
wedding budget calculator	wedding budget <b>wedding cost calculator</b> wedding calculator wedding cost breakdown <b>wedding budget worksheet</b>	monthly budget calculator online budget calculator <b>budget wedding</b> budget calculator free average wedding budget	<b>wedding planning checklist</b> wedding budget template <b>wedding budget worksheet</b> wedding checklist printable wedding costs average	<b>wedding budget worksheet</b> wedding vendors <b>the knot</b> wedding planning checklist <b>wedding wire</b>
gmat prep classes	gmat prep gmat classes gmat <b>kaplan gmat course</b> <b>kaplan gmat</b>	gmat prep online gmat prep courses <b>gmat online prep</b> online gmat prep best gmat prep courses	<b>gmat preparation courses</b> sample gmat tests <b>gmat prep class</b> which is easier gre or gmat how much is the gmat	<b>gmat study books</b> <b>gmat test prep classes</b> <b>gmat prep class</b> <b>kaplan gre courses</b> free gmat sample tests
how to build a fence	building a fence build your own fence how to build a wood fence do it yourself fence how to build a privacy fence	how to build fence build fence build a fence how to build a cheap fence build your own fence	<b>how to build a fence on a hill</b> how to fence a yard how to build a fence gate how to build a brick wall fence <b>how to build a fence video</b>	how to build a fence minecraft fancy fences and gates <b>how to build a metal fence</b> <b>home depot com fencing</b> <b>back yard fences</b>
solar panels	<b>solar panels for homes</b> <b>solar electric panels</b> <b>solar power</b> ebay solar panels how to make solar panels	solar panels for your home solar panels for residential homes <b>solar panels for</b> solar panels on ebay davis solar	<b>solar power</b> <b>solar energy</b> <b>solar panels for homes</b> <b>solar power systems</b> solar panel	<b>solar power</b> <b>solar panels for homes</b> <b>solar panels on sale</b> <b>solar panel kits</b> <b>solar panels for sale</b>

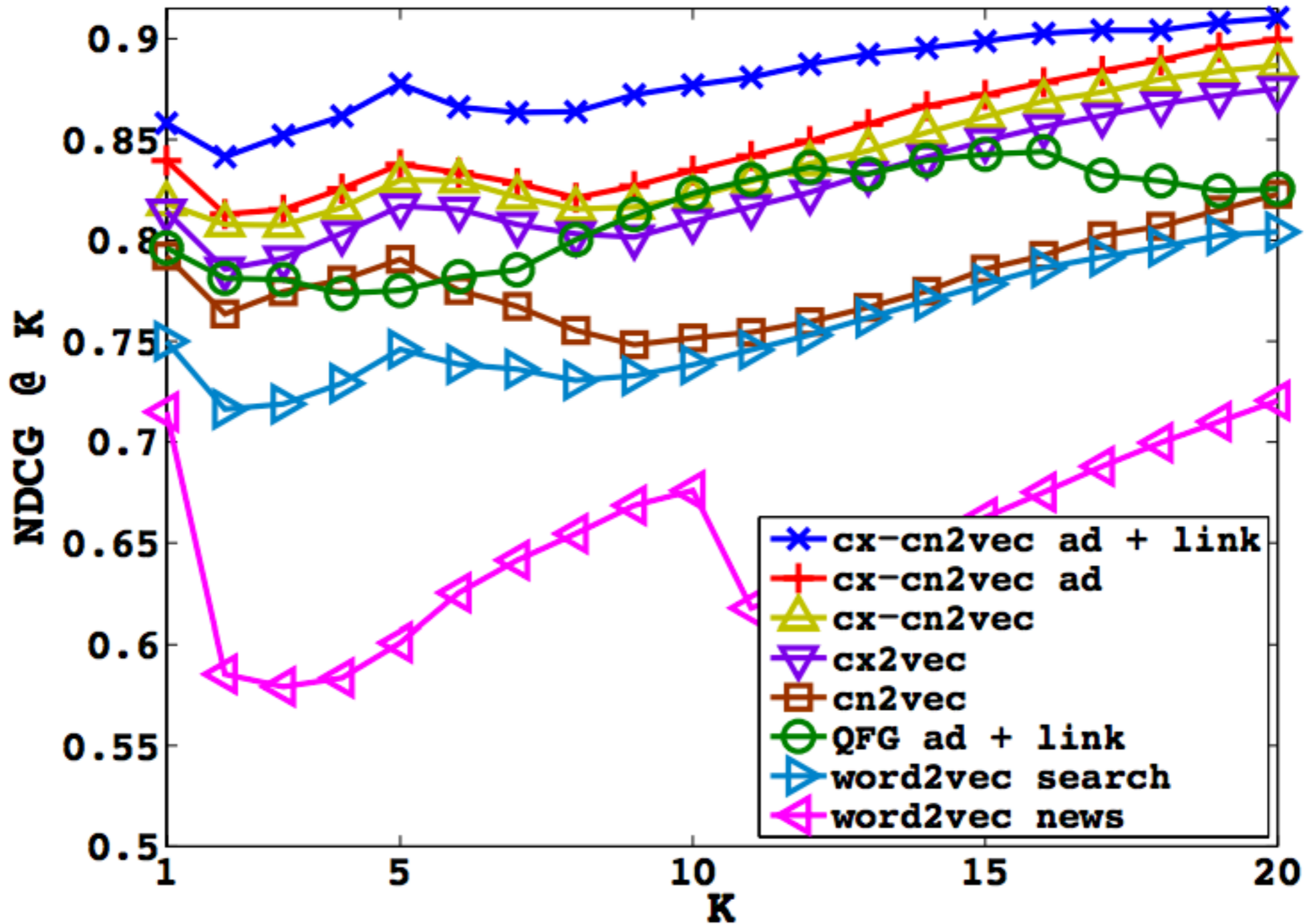
# Editorial Evaluation

## In-dictionary

grade	pairs	cn2vec	cx2vec	cx-cn2vec	cx-cn2vec <sub>ad</sub>	<b>cx-cn2vec<sub>ad+link</sub></b>	word2vec <sub>news</sub>	word2vec <sub>search</sub>	QFG <sub>ad+link</sub>
<b>Excellent</b>	1,518	0.630 (0.136)	0.658 (0.107)	0.669 (0.107)	0.668 (0.114)	<b>0.733 (0.094)</b>	0.818 (0.146)	0.648 (0.151)	0.329 (0.667)
<b>Good</b>	5,531	0.599 (0.136)	0.621 (0.125)	0.637 (0.100)	0.632 (0.104)	<b>0.683 (0.097)</b>	0.770 (0.152)	0.614 (0.155)	0.205 (0.574)
<b>Fair</b>	4,021	0.550 (0.167)	0.565 (0.129)	0.577 (0.124)	0.566 (0.130)	<b>0.605 (0.132)</b>	0.749 (0.190)	0.567 (0.173)	0.114 (0.366)
<b>Bad</b>	4,229	0.398 (0.196)	0.363 (0.170)	0.349 (0.184)	0.336 (0.187)	<b>0.425 (0.179)</b>	0.517 (0.280)	0.395 (0.201)	0.166 (0.584)
avg. p-value	-	1.39e-15	5.44e-26	8.27e-28	2.99e-30	<b>1e-100</b>	4.121e-07	1.014e-14	0.013
<b>Excellent</b>	2,119	0.791 (0.166)	0.623 (0.141)	0.628 (0.134)	0.623 (0.143)	<b>0.668 (0.147)</b>	0.824 (0.145)	0.790 (0.157)	-
<b>Good</b>	11,305	0.752 (0.155)	0.587 (0.135)	0.592 (0.130)	0.584 (0.137)	<b>0.612 (0.141)</b>	0.796 (0.145)	0.756 (0.156)	-
<b>Fair</b>	11,146	0.715 (0.136)	0.561 (0.145)	0.565 (0.139)	0.558 (0.146)	<b>0.584 (0.147)</b>	0.769 (0.175)	0.707 (0.141)	-
<b>Bad</b>	7,849	0.635 (0.199)	0.383 (0.211)	0.387 (0.209)	0.382 (0.212)	<b>0.410 (0.208)</b>	0.509 (0.311)	0.602 (0.208)	-
avg. p-value	-	4.99e-26	1.137e-27	4.0423e-32	1.196e-30	<b>2.926e-40</b>	2.038e-16	9.388e-22	-

## Out-of-dictionary

# NDCG@K Values



# Experiments on TREC Data

Method	Editorial grade	Levenshtein dist.
$QFG_{ad+link}$	1.0441	11.70
$word2vec_{news}$	0.9189	10.91
$word2vec_{search}$	0.9492	11.32
$cn2vec$	0.9571	11.37
$cx2vec$	1.1273	13.79
$cx-cn2vec$	1.1343	13.13
$cx-cn2vec_{ad}$	1.2281	13.62
<b><math>cx-cn2vec_{ad+link}</math></b>	<b>1.2457</b>	<b>13.25</b>

# Coverage and Estimated Improvements

Method	In-house data			TREC data		
	Coverage	Revenue potential	eCPM	Coverage	Revenue potential	eCPM
QFG <sub>ad+link</sub>	1.00	1.00	1.00	1.00	1.00	1.00
word2vec <sub>news</sub>	0.76	0.39	0.46	0.32	0.66	0.73
word2vec <sub>search</sub>	0.87	0.58	0.77	0.57	0.84	0.75
cn2vec	0.89	0.62	0.84	0.59	0.84	0.74
cx2vec	1.16	1.80	1.41	1.41	1.16	1.20
cx-cn2vec	1.18	1.86	1.38	1.44	1.21	1.19
<b>cx-cn2vec<sub>ad</sub></b>	<b>1.20</b>	<b>1.89</b>	<b>1.60</b>	<b>1.52</b>	<b>1.35</b>	<b>1.31</b>
cx-cn2vec <sub>ad+link</sub>	1.18	1.88	1.45	1.50	1.28	1.22





# Conclusions and Future Work

- We presented the first query rewriting mechanism using word embeddings
  - We evaluated the proposed methods using both in-house and publicly available TREC data sets
  - When compared to the current state-of-the-art approaches, we showed that context-content2vec generates the most relevant query rewrites, while at the same time maintains high level of ad coverage
- Future:
  - Better sessionization algorithms
  - Navigational query detection to bias our method less towards generating those queries





*Gracias*

