



Applying Topic Segmentation to Document-Level IR

Software Engineering
Conference Russia 2018
October 12-13
Moscow

IRELA
Gennady Shtekh
Polina Kazakova
Nikita Nikitinsky

MSU
Nikolay Skachkov

What is IR

Information Retrieval:
matching a query with
relevant documents



how to catch a unicorn

Web Images Video News More Anytime

Also try: how to catch a unicorn in sims 3

How To Catch A Unicorn - Video Results

Thumbnail	Title	Source	Duration
	HOW2: How to Catch a Unicorn!	youtube.com	2:56
	How to Catch a Unicorn	vimeo.com	1:03
	How to Catch a Unicorn (Original)	youtube.com	2:52
	How to Catch A Unicorn!	youtube.com	11:28

How To Catch A Unicorn videos

How to Catch a Unicorn - wikiHow! - The How-To Manual YOU Can Laugh At
[wikihow.com/catch-a-unicorn.html](https://www.wikihow.com/catch-a-unicorn.html)

to try to catch the unicorn at 2:09pm-3:06pm. 4. If you follow all the steps just right the unicorn should come right over and lay its head on the maidens lap. 5. If you want to take the unicorn home, alive, catch in a net and put in a horse or bear cage. 6.

How to Catch A Unicorn! - YouTube
[watch?v=ilzIHeNshe0](https://www.youtube.com/watch?v=ilzIHeNshe0)

Here is a video that I know you guys will love about Unicorns! If you read this comment what your unicorn should be with the unicorn emoji at the end. Muc...

Our Case



01 /

Document-Level
Information Retrieval:
query is also a
document

02 /

Non-conventionalized
term

03 /

Querying by
example

Our Hypothesis



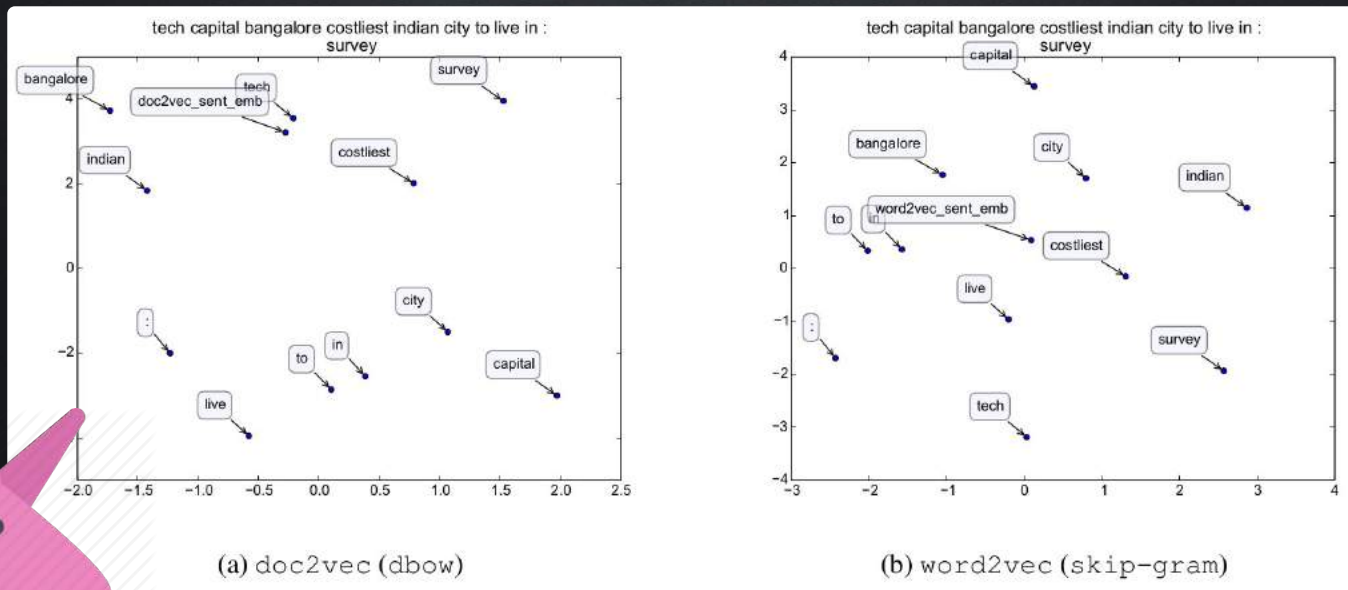
**Information
Retrieval quality
increases by using
topic segmentation
of documents.**

**Topic segmentation:
splitting texts into
semantically
homogeneous
blocks.**

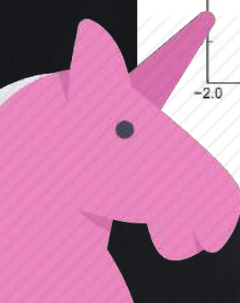
**Why to do this?
Quality of machine
learning algorithms on
short (and topically
coherent) texts is
supposed to be better.**

Why Do we Think so?

Example:
document
embeddings



Lau, Jey Han, and Timothy Baldwin.
"An empirical evaluation of doc2vec with practical insights into document embedding generation." arXiv preprint arXiv:1607.05368 (2016).



Topic Models

$$p(w | d) = \sum_{t \in T} p(w | t, d) p(t | d) = \sum_{t \in T} p(w | t) p(t | d) = \sum_{t \in T} \varphi_{wt} \theta_{td}.$$

01 /

Soft-clustering

02 /

Every document is described as a mixture of topics

03 /

Every topic is described as a mixture of words

04 /

EM-algorithm

05 /

PLSA, LDA

06 /

ARTM - regularizers

Topic Segmentation Pipeline

01 / Skachkov, N., Vorontsov, K.
Improving topic models with
segmental structure of texts.
[1]

02 / Based on ARTM
**Additive Regularization
of Topic Models [2]**
(BigARTM tool [3])

Topic Segmentation Pipeline

Part #1 /

Constructing
topic model
under sparsity
assumption

Gradual estimation of
segments borders:

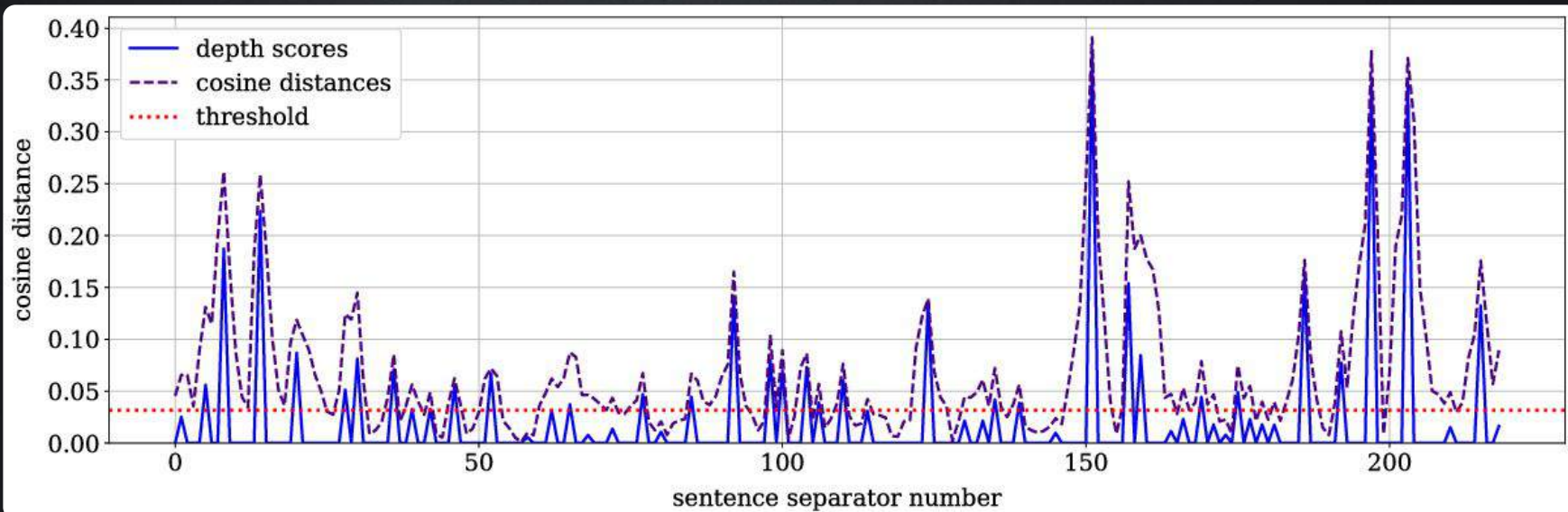
- ★ first, use sentence borders
- ★ then merge adjacent segments if they have the same topics

Topic Segmentation Pipeline

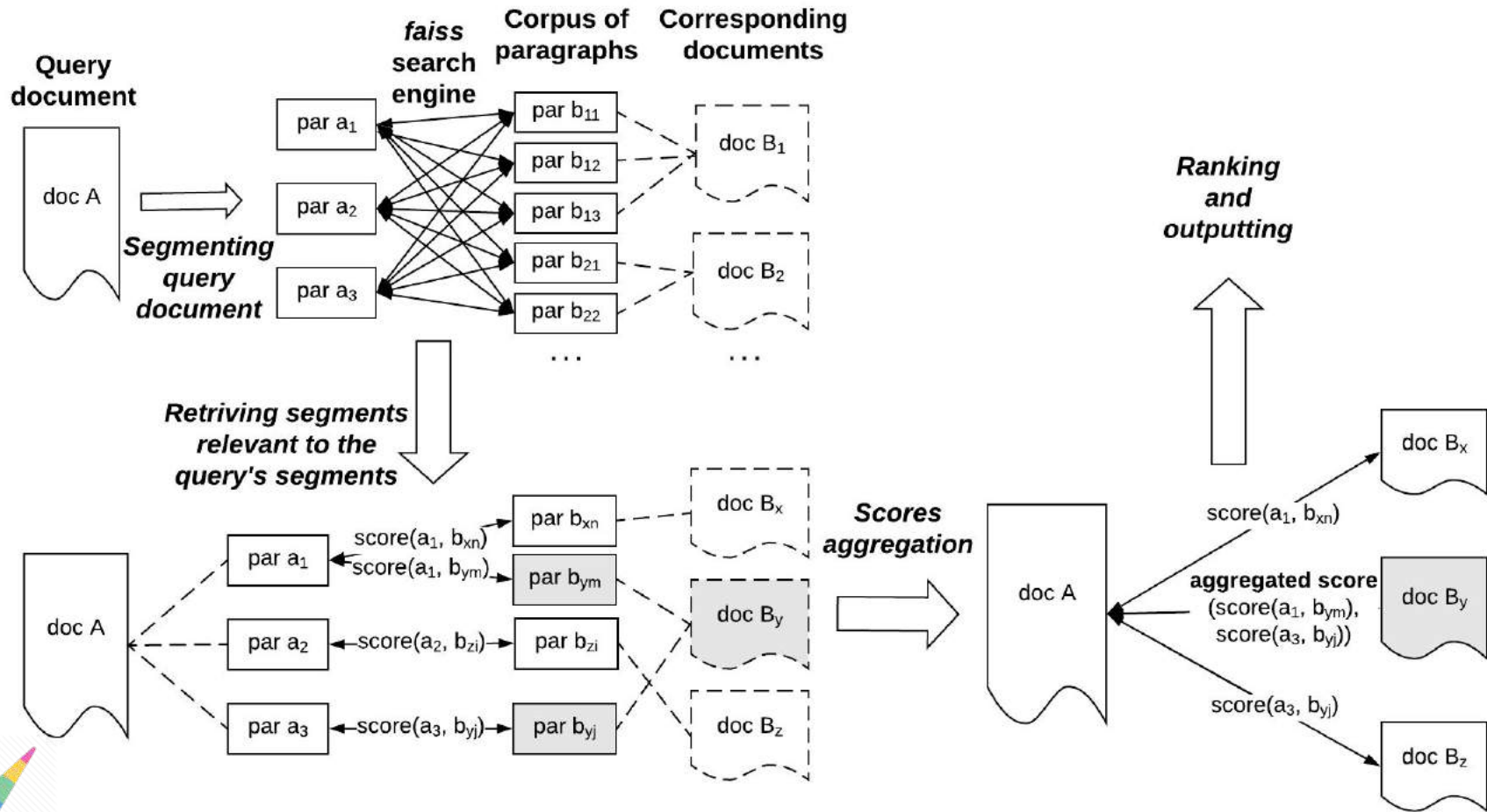
Part #2 / Topic Tiling [4] algorithm:

- ★ For each sentence boundary, consider left and right windows of a length n and compute distance
- ★ Smooth the distances - depth scores (ds)
- ★ Calculate threshold:
 - $threshold = mean(ds) - alpha * sqrt(sd(ds))$
 - $alpha$ is varied to change the granularity (default - 0.5)
- ★ Sentence separators with scores more than the threshold - segments boundaries

Topic Segmentation Pipeline



Retrieval Pipeline



Experiments: Data

Data

arXiv preprints
(140000 preprints:
train - 95000,
test - 45000)

Preprocessing

spaCy + some
manual rules
(removing
mathematical
symbols and short
strings)

Test set

triplets:
query paper -
relevant paper -
non-relevant paper;
based on arXiv subjects
(15715 triplets)

! The same technique was used in Andrew M Dai, Christopher Olah, and Quoc V Le. 2015.
Document embedding with paragraph vectors. [6]

Experiments: Baselines

- ★ ARTM segment model
(the only model trained on our data!)
- Pretrained models:
- ★ averaged word2vec
(simple and normalized)
 - ★ averaged GloVe
(simple and normalized)
 - ★ fastText
(averaged and original)
 - ★ doc2vec
 - ★ sent2vec

Experiments: Evaluation

Aggregation

- ★ **Mean:** relevance score = mean of paragraphs scores
- ★ **Best N:** relevance score = mean of N most relevant paragraphs scores (N = 1, 3, 5)

Granularity

- ★ **Two 'granulites':** coarser (alpha = 0.3) and finer (alpha = 0.5)

Evaluation

- ★ **Accuracy:** proportion of correctly ranked pairs of documents in total number of triplets

Results



<i>Model</i>	<i>No segm.</i>	<i>Finer segmentation</i>				<i>Coarser segmentation</i>			
		<i>best 1</i>	<i>best 3</i>	<i>best 5</i>	<i>mean</i>	<i>best 1</i>	<i>best 3</i>	<i>best 5</i>	<i>mean</i>
<i>ARTM</i>	0.817	0.761	0.771	0.773	0.780	0.765	0.772	0.774	0.783
<i>sent2vec</i>	0.770	0.807	0.808	0.807	0.783	0.808	0.809	0.807	0.775
<i>fastText</i>	0.751	0.784	0.785	0.782	0.684	0.784	0.785	0.782	0.680
<i>doc2vec</i>	0.814	0.783	0.785	0.782	0.628	0.780	0.781	0.778	0.636
<i>avW2V</i>	0.817	0.820	0.824	0.822	0.774	0.821	0.823	0.822	0.768
<i>avnormW2V</i>	0.580	0.584	0.583	0.587	0.620	0.584	0.582	0.586	0.616
<i>avGloVe</i>	0.779	0.779	0.779	0.778	0.712	0.777	0.778	0.777	0.709
<i>avnormGloVe</i>	0.573	0.601	0.609	0.609	0.588	0.602	0.609	0.610	0.589
<i>avfasText</i>	0.662	0.746	0.751	0.746	0.638	0.746	0.750	0.745	0.632 ₁₇

Results and observations

★ Segmentation does improve the majority of the models

★ ARTM-based model works better on whole texts

★ Small values of N in aggregation are probably better

★ Influence of segmentation granularity is unclear

★ Influence of text style perhaps must be taken into account

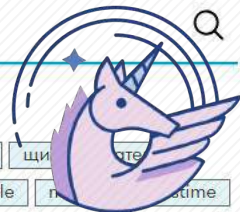
★ word2vec beats everything (???)

Real Applications: Our Experience

★ **Cross-lingual
search
engine**

★ **Ad-hoc
retrieval task**

как поймать единорога



Контекстное расширение поискового запроса

- выследить
- убить
- убежать
- перехитрить
- украсть
- достать
- заманить
- затащить
- поймал
- утащить
- дракон
- щипать
- восставать
- грифон
- олень
- леопард
- червлёный
- держать
- журавль
- коронованный
- unawares
- blacklist
- rule
- punish
- distract
- intercept
- reason
- unsuspecting
- dragon
- wyvern
- cockatrice
- basilisk
- beast
- seahorse
- dragonslayer
- bestiary
- dragonette
- minotaur

274 результата (показано 274)

рус англ

Единорог



Единоро́г или инро́г (др.-евр. דַּחַר; др.-греч. μόνό-κερως; лат. unicornis) — мифическое существо, символизирует целомудрие, в широком смысле духовную чистоту и искания. Чаще всего его представляют в виде коня с одним рогом, выходящим из лба. == Описание и символика == === У древних авторов === Самым ранним изображениям единорогов больше 4 тысяч лет, найдены в Индии. Затем стали появляться в мифах Западной Азии. В древней Греции и Древнем Риме считались реально существующими животными. Изображения единорога, попадающиеся на древнеегипетских памятниках и на скалах южной Африки, являются ри

66.959% соответствие запросу

Tue Jun 26 16:09:15 UTC 2018

Подробнее →

Unicorn



The unicorn is a legendary creature that has been described since antiquity as a beast with a single large, pointed, spiraling horn projecting from its forehead. The unicorn was depicted in ancient seals of the Indus Valley Civilization and was mentioned by the ancient Greeks in accounts of natural history by various writers, including Ctesias, Strabo, Pliny the Younger, and Aelian. The Bible also describes an animal, the re'em, which some versions translate as unicorn. In European folklore, the unicorn is often depicted as a white horse-like or goat-like animal with a long horn and cloven hoov

66.688% соответствие запросу

Sun Jul 15 15:53:12 UTC 2018

Подробнее →

References

- 01 / Nikolay Skachkov and Konstantin Vorontsov. 2018. *Improving topic models with segmental structure of texts*. In *Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference “Dialogue” (2018)*. 652–661.
- 02 / Konstantin Vorontsov and Anna Potapenko. 2015. *Additive regularization of topic models*. *Machine Learning* 101, 1-3 (2015), 303–323.
- 03 / Konstantin Vorontsov, Oleksandr Frei, Murat Apishev, Peter Romov, and Marina Dudarenko. 2015. *Bigartm: Open source library for regularized multimodal topic modeling of large collections*. In *International Conference on Analysis of Images, Social Networks and Texts*. Springer, 370–381.
- 04 / Martin Riedl and Chris Biemann. 2012. *Text segmentation with topic models*. *Journal for Language Technology and Computational Linguistics* 27, 1 (2012), 47–69.
- 05 / Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. *Billion-scale similarity search with GPUs*. *arXiv preprint arXiv:1702.08734 (2017)*.
- 06 / Andrew M Dai, Christopher Olah, and Quoc V Le. 2015. *Document embedding with paragraph vectors*. *arXiv preprint arXiv:1507.07998 (2015)*

Links



- spaCy:** spacy.io/
- word2vec:** github.com/jhlau/doc2vec#pre-trained-word2vec-models
- GloVe:** nlp.stanford.edu/projects/glove/
- fastText:** fasttext.cc/docs/en/english-vectors.html
- doc2vec:** github.com/jhlau/doc2vec#pre-trained-doc2vec-models
- sent2vec:** github.com/epfml/sent2vec#downloading-pre-trained-models

Contact us

Follow IRELA on:

Telegram: t.me/irelaru

Medium: medium.com/@irela

Facebook: facebook.com/irelaru

Gmail: kazakova1537@gmail.com

Telegram: t.me/brnzz

