

# More Than Words: Collocation Retokenization for Latent Dirichlet Allocation Models

**Jin Cheevaprawatdomrong**  
Chulalongkorn University  
jin236248@gmail.com

**Alexandra Schofield**  
Harvey Mudd College  
xanda@cs.hmc.edu

**Attapol T. Rutherford\***  
Chulalongkorn University  
attapol.t@chula.ac.th

## Abstract

Traditionally, Latent Dirichlet Allocation (LDA) ingests words in a collection of documents to discover their latent topics using word-document co-occurrences. Previous studies show that representing bigrams collocations in the input can improve topic coherence in English. However, it is unclear how to achieve the best results for languages without marked word boundaries such as Chinese and Thai. Here, we explore the use of retokenization based on chi-squared measures,  $t$ -statistics, and raw frequency to merge frequent token ngrams into collocations when preparing input to the LDA model. Based on the goodness of fit and the coherence metric, we show that topics trained with merged tokens result in topic keys that are clearer, more coherent, and more effective at distinguishing topics than those of unmerged models.

## 1 Introduction

Latent Dirichlet allocation (LDA) models provide useful insights into themes and trends in a large text collection through the unsupervised inference of *topics*, or probability distributions over unigram word types in the corpus (Blei et al., 2003). Topics from these models are often interpreted based on their highest-probability words, with documents expressed as vectors of proportions of each topic. Unfortunately, the context in which these tokens arise can be obscured in the bag-of-words rendering of text as unigram counts in documents. For instance, a topic with high probabilities of both “coffee” and “table” is tempting to interpret as focusing on the furniture item “coffee table”, but both words could be frequent in a discussion of cafes containing no coffee tables. This problem is amplified in languages without marked word boundaries, such as Chinese and Thai: while existing tokenizers in these languages can segment characters into

words, there is always a question about to what extent the tokenizers should group words together. Words that have been segmented by tokenizers may not express the concept of the original text if they were found as parts of collocations. Meaningful interpretation of topics can be lost without careful recombination of these words.

We hypothesize that the morphology of the language should play an important role in determining the suitable pre-processing steps that would improve the results of topic models. The main morphological types we consider are *synthetic language* and *analytic language*. Synthetic languages use many morphemes to compose a word and can be further divided into fusional and agglutinative languages. Fusional languages such as German differ from agglutinative languages such as Korean and Japanese: a single morpheme in fusional languages can code for many morphosyntactic features. On the other hand, analytic languages such as Thai and Chinese convey meanings by relating many words together, and morphological devices are more rarely used. Under our hypothesis, analytic languages should benefit from token merging, but synthetic languages might not because the meaning is conveyed by inflection (through bound morphemes) and agglutination (through free morphemes).

In this project, we investigate the effects of token merging as a pre-processing step, and study how those effects vary based on the writing systems and the morphological features of the languages. We evaluate three measures to determine when to merge multiple adjacent words into conceptually-unified phrasal tokens prior to LDA model training: chi-squared statistics,  $t$ -statistics, and raw frequency counts of phrases. We test these merging strategies on English, German, Chinese, Japanese, Korean, Thai, and Arabic. This set of languages is drawn from various writing systems and different morphological typology to see which type of

\*Corresponding author

language favors which type of merging strategy.

The main contributions of this paper are as follows:

- We determine through empirical studies that a  $t$ -statistic and raw-frequency approach to token merging improves the topic modeling results across all language types and writing systems for the corpora that do not differ much from the collocation training data.
- We also show the positive consequences of token merging: the percentage of merged tokens in the LDA training data is correlated with the quality of the topic modeling results.
- Finally, we provide evidence that the popular approach of applying a  $\chi^2$  measure to token merging tends to overfit to the collocation training data and result in a low percentage of merged tokens in a number of languages, making it a less suitable general-purpose approach than  $t$ -statistics.

## 2 Related Work

Pre-processing steps can substantially alter the results of the LDA models even in languages with good tokenization heuristics such as English (Schofield and Mimno, 2016; May et al., 2016). We believe that languages that do not have clear tokenization standards deserve investigation into what kind of processing is appropriate. Many works recognize that LDA results can be improved when input are including phrases (Lindsey et al., 2012; Lau et al., 2013; Yu et al., 2013; El-Kishky et al., 2014; Wang et al., 2016; Bin et al., 2018; Li et al., 2018). We consider it valuable to specifically assess approaches to determining these phrases.

Despite their popularity in analyzing large amounts of text data, LDA models are notoriously complex to evaluate. One must evaluate both the statistical fit of a model and the human-registered thematic coherence of the words found to arise in the high-probability words, or *keys*, of a topic, which may not correlate (Chang et al., 2009). Analyses often combine evaluations of fit (Wallach et al., 2009) and automated approximations of human judgments of coherence (Bouma, 2009; Mimno et al., 2011) based on mutual information, even with the expectation these may only somewhat correlate with true human judgments (Lau et al., 2014). A limitation of these existing approaches, however,

is that they expect the vocabulary and tokenization to remain constant between the two models. For our evaluation, we use a normalized log-likelihood approach to capture fit while accounting for changes in vocabulary (Schofield and Mimno, 2016).

## 3 Collocations as LDA Token

Collocations consist of two or more words that express conventional meaning, which can convey information about multi-word entities, context, and word usage. We hypothesize that the introduction of multi-word tokens, which capture collocations as bigrams or trigrams by way of concatenation of adjacent tokens, can help achieve more useful and coherent topic models. For languages without clear word boundaries, there is a possible additional benefit to multi-word tokens: it can be hard to intuit whether inferred word boundaries will have a large impact on the final results. Merging adjacent words into ‘multi-word’ tokens may help remedy the potential problem of a segmentation that is not optimal for topic modeling purposes.

Many methods are possible to select collocations to merge from tokenized text (Manning and Schutze, 1999). In this paper, we evaluate the chi-squared statistics ( $\chi^2$ ), the  $t$ -statistic and raw frequency as approaches to develop a threshold for merging collocations into multi-word tokens prior to topic model training. The chi-squared measure  $\chi^2(w_1, w_2)$  and  $t(w_1, w_2)$   $t$ -statistic for two adjacent tokens  $w_1$  and  $w_2$  are defined as:

$$\chi^2(w_1, w_2) = \frac{(P(w_1, w_2) - P(w_1)P(w_2))^2}{P(w_1)P(w_2)} \quad (1)$$

$$\begin{aligned} t(w_1, w_2) &= \frac{\bar{x} - \mu}{\frac{s^2}{N}} \\ &\approx \frac{P(w_1, w_2) - P(w_1)P(w_2)}{\sqrt{\frac{P(w_1, w_2)}{N}}} \quad (2) \end{aligned}$$

We first compute the collocation measures for all bigrams on a large collocation training corpus. Then we select the top bigrams that score the highest on the collocation measures and add those to our lexicon. After we tokenize and pre-process the collection of documents on which we would like to train LDA, we retokenize the data based on the collocation training corpus. We find all of the bigrams in the LDA training data that are also found in the top bigram lexicons that we obtain from the

collocation training corpus. Then, the LDA training process proceeds as usual but with some of the original tokens merged into multi-word tokens as defined from the collocation training data.

#### 4 Evaluation Metrics

We consider two primary evaluation metrics for exploring the effect of merging tokens: one based on log-likelihood, and one based on silhouette coefficients.

**Held-Out Likelihood.** When multi-word phrases are converted to individual tokens, the number of tokens in the document decreases while the size of the corpus vocabulary increases. It is therefore illogical to compare the likelihoods of the word-token model and collocation-token model directly. In order to normalize the scores between the two models that do not have the exact same vocabulary and tokens, we use the log-likelihood ratio between the LDA model likelihood and the null (unigram) likelihood for each model. In other words, we normalize the LDA model likelihood ( $\mathcal{L}_{\text{model}}$ ) by dividing it with the unigram likelihood ( $\mathcal{L}_{\text{unigram}}$ ) as introduced by Schofield and Mimno (2016). Therefore, the normalized loglikelihood per token (PTLL<sub>norm</sub>) is

$$\text{PTLL}_{\text{norm}} = \frac{\log \mathcal{L}_{\text{model}} - \log \mathcal{L}_{\text{unigram}}}{N} \quad (3)$$

where  $N$  is the number of tokens. Since likelihood per token has been normalized by the unigram likelihood per token, the higher the PTLL, the better the model.

**Concatenation-based Embedding Silhouette (CBES)** Previous measures of topic coherence rely on statistics from the training data and assume that the vocabularies are identical for both models, which is not the case for our settings. To address this, we propose a new application of the silhouette coefficients (Rousseeuw, 1987), a common clustering evaluation metric to measure topic coherence.

A good topic should have all of its topic keys close to each other and away from other words that do not belong in the same topic. Therefore, the word embeddings of these topic keys should have shorter cosine distances within the same topic, and longer distances to the topic keys in other topics. When words are represented as a vector, this is exactly what the silhouette coefficients measure. To compute them, we first compute the  $a(i)$ , which is the mean cosine distance between topic-key  $i$

and other topic-keys in the same topic.

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, i \neq j} d(i, j) \quad (4)$$

where  $d(i, j)$  is the distance between  $i$ th and  $j$ th topic-key and  $|C_i|$  is the number of topic-keys in topic  $i$ . Then for each other topic, we compute the mean of the distance of topic-key  $i$  to topic-keys in that other topic. And  $b(i)$  is the smallest of such mean among other topics.

$$b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j) \quad (5)$$

After obtaining  $a(i)$  and  $b(i)$ , the silhouette coefficient for topic-key  $i$  is defined as:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}, \text{ if } |C_i| > 1 \quad (6)$$

and

$$s(i) = 0, \text{ if } |C_i| = 1 \quad (7)$$

The silhouette coefficient for the entire model is the average  $s(i)$  over all  $i$ . The larger silhouette coefficient means that topic-keys are relatively similar within their topic and different from other topics.

In order to compare the distances among words merged by different criteria, all compared word embeddings must be in the same space. Since merged tokens will modify the vocabulary of the corpus, we create four versions of the word embedding training corpus: the original version and the three other versions where tokens are merged based on  $\chi^2$ ,  $t$  and frequency collocation measures. We train the word embeddings on these four versions of the corpus so we can then compare word embeddings on a consistent vocabulary in each retokenization scheme.

#### 5 Experiments

We hypothesize that morphology should play an important role in determining the suitable preprocessing steps. We test our methods on one fusional language (German), two agglutinative languages (Japanese and Korean), three analytic languages (Chinese, Thai, and Arabic), and English, which can be thought of as either analytic or fusional. These languages also represent languages drawn from all writing systems: logograms (Chinese), syllabic system (Japanese), featural system (Korean), abugida (Thai), abjad (Arabic), and true alphabets (English and German).

	Domains	Docs (K)	Tokens (M)	%Merged		
				CHI	T	FREQ
EN-NYTimes	News	53	0.7	1.64	12.71	12.72
EN-SOTU	Speeches	42	0.8	0.86	9.76	10.33
EN-Yelp	Restaurants	67	2.1	0.16	7.85	8.97
DE-10kGNAD	News	222	1.9	0.09	7.46	7.68
CN-Chinanews	News	49	0.8	0.00	11.61	11.64
CN-Dianping	Restaurants	40	0.8	0.01	2.82	2.80
CN-Douban	Movies	98	0.6	0.03	4.17	4.23
JA-JapanNews	News	528	3.6	21.74	21.95	21.85
KO-KAIST	Misc	20	0.2	19.82	20.71	21.27
TH-Prachathai	News	32	4.4	0.07	15.97	14.06
TH-Wongnai	Restaurants	40	1.2	0.00	8.52	6.09
TH-BEST	Misc	7	2.1	0.03	14.94	13.09
TH-TNC	Misc	4	1.0	0.03	13.65	12.00
AR-ANT	News	60	1.1	0.16	26.13	27.45

Table 1: A survey of corpora providing the number of documents and tokens, as well as the percentage of unigram tokens merged using each approach.

The English corpora are drawn from The New York Times (Sandhaus, 2008), the Yelp Dataset<sup>1</sup>, and United States State of the Union addresses (1790 to 2018) divided into paragraphs<sup>2</sup>. The German data come from Ten Thousand German News Articles Dataset<sup>3</sup>. The Chinese data come from three corpora: the news articles from Chinanews<sup>4</sup>, restaurant reviews from Dianping<sup>5</sup>, and the movie reviews from Douban<sup>6</sup>. The Japanese data is from the Webhose’s Free Datasets<sup>7</sup>. The Korean data come from the KAIST Corpus<sup>8</sup>. The Thai data come from the news articles in Prachathai<sup>9</sup>, the restaurant reviews from Wongnai<sup>10</sup>, the BEST corpus<sup>11</sup>, and the Thai National Corpus (Aroonmanakun, 2007). The Arabic data come from the Antcorpus (Chouigui et al., 2017). Each corpus is separated into 75% training documents and 25% test documents (Table 1).

We train the  $\chi^2$ ,  $t$ , and frequency-based tokenizers for each language on Wikipedia articles for that language. For all languages, we use the reduced version of Wikipedia database, except for English we use the filtered Wiki103 dataset (Merity et al., 2016). English, German, Chinese, Japanese, Korean, Thai and Arabic documents are tokenized with NLTK (Bird, 2006), SoMaJo (Proisl and

<sup>1</sup>www.yelp.com/dataset

<sup>2</sup>www.kaggle.com/rtatman/state-of-the-union-corpus-1989-2017

<sup>3</sup>github.com/tblock/10kGNAD

<sup>4</sup>www.chinanews.com

<sup>5</sup>github.com/zhangxiangxiao/glyph

<sup>6</sup>www.kaggle.com/utmhikari/doubanmovieshortcomments

<sup>7</sup>webhose.io/free-datasets/japanese-news-articles/

<sup>8</sup>semanticweb.kaist.ac.kr/home/index.php/KAIST\_Corpus

<sup>9</sup>github.com/PyThaiNLP/prachathai-67k

<sup>10</sup>www.kaggle.com/c/wongnai-challenge-review-rating-prediction

<sup>11</sup>thailang.nectec.or.th/downloadcenter

	$\chi^2$ - $t$	$\chi^2$ -freq	$t$ -freq
English	8.90	7.78	74.87
German	0.00	0.00	83.06
Chinese	0.00	0.00	86.48
Japanese	29.06	22.60	73.34
Korean	10.56	7.34	71.95
Thai	0.22	0.06	67.25
Arabic	1.22	1.20	66.89

Table 2: The percentage of overlapping merged tokens between two methods of retokenization computed on the retokenization training data.  $t$  and  $\chi^2$  yield similar results for all languages.

Uhlig, 2016), Stanford Word Segmenter (Tseng et al., 2005), Fugashi (McCann, 2020), KoNLPy (Park and Cho, 2014), Attacut (Chormai et al., 2020) and Camel-tools (Obeid et al., 2020) respectively. For each criterion, we create a list of 50,000 top bigrams that have the highest scores. These lists of top bigrams will be used to merge words in the input of the LDA, effectively training a new tokenizer.

To train word embeddings, we use the gensim (Řehůřek and Sojka, 2010) implementation with the Continuous Bag-of-Word (CBOW) algorithm (Mikolov et al., 2013) to obtain word embeddings. The training corpora and their collocation versions are prepared based on the tokenizers that we discuss above. We preprocess the word embedding training data and the LDA training data the same way. For English, we lemmatize and lowercase the data. For Korean, Japanese, and Arabic, we lemmatize the data. For German, Chinese, and Thai, we do not do any normalization.

We use MALLET (McCallum, 2002) implementation of LDA with the default hyperparameters to train and evaluate topic models in both word and multi-word (collocation) documents with 10, 50, 100 topics. We run the experiment 3 times for each combination of corpus, type of retokenization (no retokenization,  $\chi^2$ ,  $t$  or frequency) and number of topics to compute the means of the normalized held-out likelihood and CBES, discussed in section 4.

## 6 Results and Discussion

The normalized log-likelihood per token of the  $t$  and frequency-based retokenization is significantly higher than the baseline for English, German, Chinese, Japanese, Korean, and Arabic for all text collections and the number of topics except EN-Yelp, TH-BEST, and TH-TNC (Table 3). Frequency-



	10 topics				50 topics				100 topics			
	Word	$\chi^2$	$t$	freq	Word	$\chi^2$	$t$	freq	Word	$\chi^2$	$t$	freq
EN-NYTimes	.3646	.3675	.4119	<b>.4386</b>	.5214	.5225	.5766	<b>.6128</b>	.5588	.5533	.6050	<b>1.0492</b>
EN-SOTU	.2699	.2660	.2967	<b>.3145</b>	.3809	.3809	.4122	<b>.4430</b>	.4135	.4101	.4367	<b>.4705</b>
EN-Yelp	.1597	.1607	.1833	<b>.2021</b>	.2589	.2599	.2893	<b>.3169</b>	.3357	.2822	.3130	<b>.3412</b>
DE-10kGNAD	.4982	.5001	.5233	<b>.5251</b>	.7272	.7272	.7622	<b>.7651</b>	.7784	.7809	.8122	<b>.8188</b>
CN-Chinanews	.5033	.5046	.5510	<b>.5592</b>	.7647	.766	.8170	<b>.8344</b>	.8427	.8394	.8847	<b>.9044</b>
CN-Dianping	.2557	.2574	.2644	<b>.2659</b>	.3899	.3906	.3965	<b>.4013</b>	.4188	.4212	.4255	<b>.4263</b>
CN-Douban	.2966	.2955	.3076	<b>.3092</b>	.4048	.4073	.4144	<b>.4173</b>	.4294	.4301	.4332	<b>.4374</b>
JA-JapanNews	.4540	<b>.7803</b>	.5942	.6342	.7173	.9268	.9339	<b>.9926</b>	.8088	1.0325	1.0316	<b>1.1003</b>
KO-KAIST	.2901	<b>1.0315</b>	.4589	.5442	.6446	.6833	.7152	<b>.8390</b>	.4755	.7437	<b>1.3443</b>	.9221
TH-Prachathai	.4367	.4331	<b>.4756</b>	.4743	.7052	<b>.8458</b>	.7699	.7719	.7854	.7854	.8537	<b>.8548</b>
TH-Wongnai	.2048	.2013	<b>.2225</b>	.2192	.3237	.3222	<b>.3472</b>	.3399	.3467	.3463	<b>.3720</b>	.3636
TH-BEST	<b>.6995</b>	<b>.6995</b>	.6704	.6838	.9148	.9190	.9279	<b>.9389</b>	.9812	.9819	.9967	<b>1.0100</b>
TH-TNC	.7420	<b>.7422</b>	.7079	.7239	.9969	.9952	1.0079	<b>1.0219</b>	1.0508	1.0473	1.0608	<b>1.0758</b>
AR-ArabNews	.3183	.3152	.4676	<b>.5663</b>	.4923	.4913	.7175	<b>.8742</b>	.5417	.5409	.7681	<b>.9355</b>

	10 topics				50 topics				100 topics			
	Word	$\chi^2$	$t$	freq	Word	$\chi^2$	$t$	freq	Word	$\chi^2$	$t$	freq
EN-NYTimes	.0143	.0153	.0246	<b>.0453</b>	-.0582	-.0625	-.0544	<b>-.0487</b>	-.0876	-.0875	-.0783	<b>-.0780</b>
EN-SOTU	.0034	-.0013	.0070	<b>.0100</b>	-.0602	-.0597	-.0595	<b>-.0527</b>	-.0812	-.0823	-.0793	<b>-.0743</b>
EN-Yelp	-.0634	-.0548	-.0465	<b>-.0337</b>	-.1117	-.1085	-.1023	<b>-.0952</b>	-.1299	-.1290	-.1179	<b>-.1153</b>
DE-10kGNAD	-.0209	-.0244	-.0190	<b>-.0134</b>	-.0804	-.0860	-.0785	<b>-.0680</b>	-.0753	-.0730	-.0655	<b>-.0599</b>
CN-Chinanews	.0002	.0018	.0152	<b>.0162</b>	-.0523	-.0559	-.0456	<b>-.0388</b>	-.0699	-.0712	-.0665	<b>-.0620</b>
CN-Dianping	<b>-.0708</b>	-.0854	-.0714	-.0744	<b>-.1278</b>	-.1316	-.1317	-.1339	<b>-.1373</b>	-.1439	-.1446	-.1439
CN-Douban	-.0226	-.0140	<b>-.0078</b>	-.0095	<b>-.0847</b>	-.0854	-.0864	-.0850	<b>-.1037</b>	-.1041	-.1073	-.1053
JA-JapanNews	-.0925	-.0655	-.0562	<b>-.0133</b>	-.1503	-.1010	-.0977	<b>-.0716</b>	-.1644	-.1120	-.1106	<b>-.0915</b>
KO-KAIST	-.0608	-.0315	-.0317	<b>-.0191</b>	-.0895	-.0691	-.0664	<b>-.0503</b>	-.0868	-.0698	-.0726	<b>-.0592</b>
TH-Prachathai	<b>-.0039</b>	-.0092	-.0040	.0160	-.0806	-.0797	-.0684	<b>-.0623</b>	-.1137	-.1121	-.0939	<b>-.0896</b>
TH-Wongnai	<b>-.0667</b>	-.0672	-.0733	-.0726	-.1468	-.1530	<b>-.1462</b>	-.1505	-.1761	<b>-.1709</b>	-.1738	-.1767
TH-BEST	-.0278	-.0187	-.0248	<b>-.0095</b>	-.0987	-.0977	-.0987	<b>-.0927</b>	-.1145	-.1153	-.1086	<b>-.1007</b>
TH-TNC	-.0284	-.0324	<b>-.0133</b>	-.0271	-.1079	-.1053	-.1332	<b>-.0964</b>	-.1281	-.1274	-.1297	<b>-.1175</b>
AR-ArabNews	-.0695	-.0673	-.0496	<b>.0124</b>	-.1255	-.1129	-.0834	<b>-.0434</b>	-.1355	-.1309	-.1010	<b>-.0735</b>

Table 3: Normalized unigram log-likelihood per token (top) and Concatenation-based Embedding Silhouette (CBES) scores (bottom) for between the baseline and retokenization models:  $\chi^2$ , texttitt, and raw frequency. Shaded cells mean that the results are inferior to the baseline, while bolded cells show the best results for each corpus and number of topics.

based retokenization gives the best results for most settings but not significantly higher than  $t$  retokenization. However, we observe mixed results from  $\chi^2$  retokenization for some languages. This is quite surprising because raw frequency was previously found to be an inferior measure of collocation. This suggests that  $t$  and frequency-based retokenization might be a more reliable method for improving the goodness of fit of the LDA model. This also suggests that Japanese and Korean might have some specific quality that interacts well with all three types of retokenization.

Similarly, we observe a general improvement in coherence for the  $t$  and frequency retokenization (Table 3). The higher CBES score indicates that topic-keys are more semantically coherent and topics are more distinct. The coherence improves after  $t$  and frequency-based retokenization for English, Japanese, Korean, and Arabic corpora regardless of the number of topics. The improvement for Thai is

spotty, and Chinanews is the only Chinese corpus in which we see improvement. This suggests that the choice of retokenization strategy might depend on the language types or the content of corpora itself. Consistent with the normalized log-likelihood results, Japanese and Korean corpora interact well with all three types of retokenization, suggesting that the morphology or typology of these two languages consistently benefit from collocation before training LDA models.

What could account for this discrepancy across languages and corpora? First, we observe a large variation of percentages of merged tokens across corpora. Because we fix the number of bigrams types to merge during the tokenizer training process to 50,000 for all three criteria (Table 1), we can use this analysis to find trends in the relative frequency of merged tokens. We see that  $\chi^2$  retokenizer only merges barely 1% of all the tokens before training the LDA models for English, Chinese,

$\chi^2$ : <b>dvenadsat apostolov, jormp jomp, malwae tweep, aboul gheit, achduth vesholom, adavari matalaku, adeste fideles, afforementionede oughtt, agoraf drws, aht urhgan, akanu ibiam, aksak maboul, alberthiene endah, alfava metraxis, alfonsas eidintana, allasani peddana, alteram partem, amantes clandestinos, amarin winitchai, amel oluna</b> <b>t</b> : <b>united states, new york, world war</b> , km h, take place, miles km, <b>los angeles, united kingdom</b> , first time, high school, tropical storm, <b>new zealand, war ii</b> , video game, mph km, h mph, <b>north america, air force</b> , two years, peak number <b>frequency</b> : <b>united states, new york, world war</b> , km h, take place, miles km, first time, <b>los angeles, united kingdom</b> , high school, tropical storm, <b>new zealand</b> , video game, <b>war ii</b> , mph km, two years, h mph, <b>north america, air force</b> , peak number
$\chi^2$ : うそ寒い 肌寒, きっこん ばったん, ざらり ぐらり, へへへへ へへへ, アウレオルス ボンバストゥス, アジ タケサカンバリン, アッシュアルク アルアウサト, アトミズム アドリアシン, アドリアシン アドリアマイシ ン, アルパイ オザラン, アワサカ ツマオ, イブリツモマブ チウキセタン, ウダヤン プラサッド, ウラマツ サ ミタロウ, エウグランディナ ロセア, エストラムスチン エストラサイト, オクタクロルテトラヒドロ メタノ フタラン, オドネ センデロル, オランバヤル ビャンバジャブ, クツミ ソクチュウ <b>t</b> : 年月, る 居る, 月日, る 事, 其の 後, 成る 居る, 昭和年, 事 出る, 年 昭和, 於 くり, 年年, 成る, 事 有る, 事 成る, 使用る, 物 有る, 存在る, 平成 年, 第 回, る 年 <b>frequency</b> : る 居る, 年月, 月日, る 事, る 年, 年年, 成る 居る, 居る 年, 其の 後, 事 有る, 昭和 年, る, る 其の, 事 成る, 事 出る, 年 昭和, 有る 年, 成る, 使用る, 於 くり
$\chi^2$ : 가닛 알훤소, 가윗일 봇일, 가즈테루 우루샤, 가톨리곤 엠블, 갈끔 가실끔, 갈라람 알부담, 감민 월민, 감성 체 널@21, 갑복 갑규, 강첸 키송, 강취완 강취일, 강홍업 강효업, 강홍선 강홍익, 개영 개영, 개초항 거륵항, 개 튀의알 똥퍼먹는, 객렬액 겁렬액, 겐러 리@KCUA, 겐런에서 겐런으로, 거대유방증 대유방 <b>t</b> : 적 인, 하다 수, 한 다, 위 한, 말 하다, 시작 하다, 사용 하다, 못 하다, 수 없다, 위치 한, 하다 않다, 사용 되다, 하 다 위해, 가지 고, 기도 하다, 일반 적, 되다 않다, 존재 하다, 기록 하다, 은 대한민국 <b>frequency</b> : 적 인, 하다 수, 하다 하다, 한 다, 사용 하다, 말 하다, 시작 하다, 하다 않다, 위 한, 못 하다, 수 없다, 위 치 한, 하다 위해, 하다는, 사용 되다, 기록 하다, 되다 않다, 하다 되다, 기도 하다, 활동 하다

Figure 1: The top 20 collocations from each retokenization methods.  $\chi^2$  favor proper names (bold-faced) more heavily than the other two methods.

German, Arabic, and Thai corpora, possibly introducing noise in the data that yield the results similar to or worse than the baseline. In contrast, the  $t$  and frequency-based retokenizers merge around 8%- 15% of all the tokens for English, German, and Chinese. Arabic has seen the highest merging percentage of 26%-27%. Notably, around 20 % of tokens are retokenized by all three retokenizers in Japanese and Korean. The truncation of the top  $\chi^2$  bigrams list might cause this different behavior. The number of  $\chi^2$  collocations that pass the hypothesis testing is significantly larger than that of  $t$  collocations. For example, there are 3.73 million  $\chi^2$  collocations versus 231 thousand  $t$  collocations in Thai for the same significance level  $\alpha = 0.005$ . This full list of  $\chi^2$  collocations includes all the top collocations from the  $t$  score and frequency treatments, implying that were we to use this significance threshold, the percentage of merged word would be at least as high as the two methods. However, the large vocabulary that the  $\chi^2$  approach induces is impractical in many applications, suggesting it is an inefficient approach if the goal is primarily to merge frequent ngrams.

Another possible effect these results may show is that the writing system or the morphology could account for this notable discrepancy in retokenization percentage across languages. For English, the top 20  $\chi^2$  collocations are primarily specific

named entities, but the  $t$  and frequency-based retokenizers yield more general compound nouns and common phrases (Figure 1). As the top 50,000  $\chi^2$  collocations contain primarily rare words, these are expected to co-occur rarely enough that even a few co-occurrences can trigger significance. Therefore, when we use this truncated list of rarely-occurring  $\chi^2$  collocations, we generally see a very low merged token percentage.

The quality of retokenization impacts both the goodness of fit the model, as indicated by the normalized log-likelihood score, and the coherence of the model, as indicated by the CBES score. Within the same language, news corpora have higher percentages of merged words when merged with  $t$  and frequency collocations, while corpora containing restaurant and movie reviews tend to see lower percentages (Table 1). This could be because the news corpora are in a similar domain to that of the Wikipedia which we use to build the list of co-occurring words. A good retokenizer (in our cases, trained on Wikipedia data) should generalize well and recognize many collocations in a new corpus, which differs somewhat from the retokenizer training data. We found a significant positive correlation between merge percentage and the margin of improvement over the baseline (the difference between the PTLL of the model without retokenization and the PTLL or CBES of the model

	Word	$\chi^2$	<i>t</i>	Freq
EN-SOTU	security social program system benefit welfare legislation need must reform propose congress health retirement administration meet enact national work insurance	health care security social insurance welfare work americans reform system cost benefit program must make need help plan pay retirement	americans <b>social_security</b> <b>health_care</b> cost families benefit pay plan save system american help reform care retirement tax medicare work coverage work make	<b>social_security</b> welfare <b>health_care</b> system benefit families insurance reform cost care health save americans retirement medicare work must pay coverage workers
DE-10kGNAD	de Spanien Madrid spanischen El Barcelona Mexiko Messi Brasilien Chile Rousseff spanische Valencia Venezuela Kuba La USA Präsidentin Real Luis	FC Der Madrid Barcelona Bayern Real Gruppe City Manchester League München Hinspiel Tore United In Die Spanien Minute Trainer Champions	Der Hinspiel Madrid Bayern Spanien Barcelona spanischen Valencia <b>Real_Madrid</b> Atletico <b>Champions_League</b> Messi Real <b>FC_Barcelona</b> Trainer Tore Saison Liverpool Gesamtscore Arsenal	Der Trainer Hinspiel Die Janko <b>Champions_League</b> Alaba Bayern Valencia Minute Saison Tore <b>Real_Madrid</b> Madrid Atletico Messi Real David Barcelona <b>FC_Barcelona</b>
CN-Chinanews	世界杯 巴西 国际 足联 时间 南非 比赛 足球 预选赛 球队 届 中新网 凌晨 抽签 进行 北京 强 小组 欧洲 支	世界杯 巴西 南非 时间 足球 欧洲 比赛 届 场 球队 德国 杯 支 球场 球迷 葡萄牙 非 洲 法国 阿根廷 凌晨	世界杯 国际 足联 巴西 足 球 南非 球迷 主席 巴西 世 界杯 体育 法国 球场 天 俄 罗斯 南非 世界杯 民众 次 卡塔尔 布拉特 标志 德国	世界杯 巴西 国际 足联 巴 西 世界杯 女足 足球 南非 国足 昨天 北京 时间 球迷 国家队 男足 中国 中国队 今 天 预选赛 小组赛 强 球队
JA-JapanNews	月 日 御 年 為 る 期 間 限 定 居 る 発 売 中 商 品 販 売 時 キ ャ ン ペ ー ン 購 入 円 成 る 下 さ る 頂 く 掲 載	億 円 大 手 事 業 商 品 社 利 益 日 社 長 企 業 市 場 同 社 米 投 資 サ ー ビ ス 同 プ ラ ン ド 中 国 因 る 海 外 販 売	月 成 る 円 店 今 年 同 億 円 年 比 前 年 増 調 査 日 本 増 え る 販 売 利 益 年 月 日 料 金 別	御 キ ャ ン ペ ー ン 商 品 ロ ッ テ 限 定 円 纏 め 掲 載 期 間 限 定 此 の 下 記 下 さ る 無 料 実 況 期 間 中 ク ー ポ ン セ ー ル 開 催 為 る 居 る 買 う 頂 く
KO-KAIST	- 교수 하다 대 연구 는 팀 한 되 다 연 황 서울대 현 광고 미 국 연구원 신 밝히다 검증 이 라고 수 첩	교수 서울대 황 대 교수 는 연구 팀 말 하다 인 도 줄기세포 취재 검증 본보다 대한 수 첩 연구원 이 르 다 논문 논란	교수 황 서울대 팀 연구 광고 대 취재 수 첩 인 줄기세포 교수 는 김 본보다 연구원 한편 프로그램 한 논문 이 라고 밝히다	교수 서울대 황 팀 연구 대 교수 는 줄기세포 연구원 수 첩 의학 취재 신 한 인 본보다 검증 논문 말 하다 대한
TH-Prachathai	แดง เลือ คน แกนนำ ชุมม นปช. สี ไทย นิโร ทษกรม เมือง เหลือ เวที รัฐบาล ทหาร สลาย จำนวน เคลื่อนไหว ปราศรัย เลือแดง ประชาชน	ชุมนุม คน เลือแดง แกนนำ เวลา นปช. ตำรวจ หน้า สลาย น. เวที สถานการณ์ พื้นที่ ประกาศ ปราศรัย เหตุการณ์ บริเวณ ประชาชน ถูกเงิน	ชุมนุม คน เลือแดง แกนนำ คน สลาย รัฐบาล เลือแดง ประชาชน เหตุการณ์ รุนแรง นปช. พื้นที่ เรียกร้อง เวที เจรจา เคลื่อนไหว ประกาศ สถานการณ์ กปปส. ตำรวจ	ชุมนุม แกนนำ คน คน เลือแดง นปช. เวที เวลา น. ปราศรัย เลือแดง รัฐบาล สลาย กปปส. หน้า เคลื่อนไหว ประกาศ เดินทาง เรียกร้อง ประชาชน ตำรวจ พันธมิตรฯ
AR-ANT	شركة في غاز طاقة أن تونسي كهرباء إنتاج من نشاط مدير كهربائي عن سبب مصنع منجم عام بتزولي متجدد	شركة في غاز طاقة من إنتاج تونسي كهرباء عن نشاط مدير مصنع أن عامل منجم صناعه بتزولي قطاع متجدد عام	شركة إنتاج تونسي نشاط غاز مؤسسة مصنع وطني عامل منجم قفص بتزولي منجم إلى الفسفاط طاقة متجدد مدير عام استغلال مائه ماء مغربي	تونسي غاز شركة كهرباء نشاط قفص منجم إنتاج مصنع إقليم سبب صباح إلى ساعة ثاني بعد رو إل بتزولي طاقة متجدد طاقة أفريل شغل الفسفاط يوم أحد سيطرة ذاتي

Figure 2: Topic keys comparison in languages.

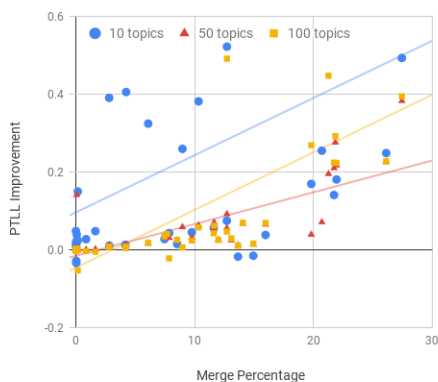


Figure 3: PTL improvement vs. merged percentage.

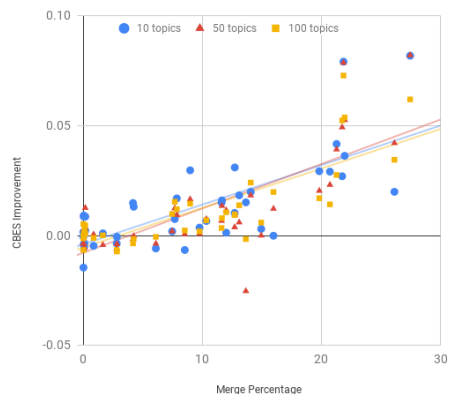


Figure 4: CBES improvement vs. merged percentage.

with retokenization). Pooling across all languages and corpora, we found the correlation coefficients of 0.41, 0.77, and 0.68 for the models with 10, 50, and 100 topics respectively for PTLT. As for the coherence metric, we found the correlation coefficients of 0.73, 0.76, and 0.79 for the models with 10, 50, and 100 topics respectively for CBES. This means the models with higher merge percentages are better than their corresponding word models in reproducing the statistics of the held-out data. This suggests that the quality of the LDA models depends on the generalizability of the retokenizers.

The LDA model results become more understandable when certain tokens are retokenized. We see merged tokens in the topic key sets of almost all topics in all corpora when retokenized based on  $t$  or raw frequency. Many of these represent non-compositional meanings that might have been lost without retokenization: for example, the collocation “social security” is not fully represented by the individual tokens “social” or “security” separately. More strikingly, the collocation ‘*kōn sūa dāng*’ refers to a political movement group in Thailand. When it is separated into *kōn* (people) *sūa* (shirt) *dāng* (red), the key meaning is totally lost. When we compare by looking at the topic-keys of the word and multi-word models, we can come up with similar topics because we as a human who understands English and has general knowledge of the world can make the connection based on surrounding topic-keys even though they are not explicitly merged. However, if we want to use these topic keys as input to other downstream tasks such as information retrieval or text classification, the merged tokens help retain the specificity of the “red shirt people” as a meaningful entity distinct from the phrase’s constituting parts.

## 7 Conclusion

In this work, we improve the quality of LDA models by better processing the input text before training the model. We found that the retokenizers trained based on  $t$  statistics and raw frequency yield an improvement across all languages considered in this study, while the  $\chi^2$  approach was a less efficient approach that focuses more on rare named entities than common noun phrases. Using retokenizers ensures that LDA models can fit better to the data, the topic keys are more coherent, and the topics are more distinct. Outputs from retokenization with  $t$  statistics and frequency approaches yield common

noun phrases in the most frequent terms of topics that represent a significant aid to both direct topic interpretation and expected utility of these topics in downstream tasks.

## Acknowledgments

This project is partially supported by Grants for Development of New Faculty Staff, Ratchadaphiseksomphot Endowment Fund. The authors would like to thank Vincent Ng, who provided us with very insightful comments as a Student Research Workshop mentor. We are also grateful for the suggestions from the anonymous reviewers from the previous submission.

## References

- Wirote Aroonmanakun. 2007. Creating the thai national corpus. *MANUSYA: Journal of Humanities*, 10(3):4–17.
- GE Bin, Chun-hui HE, Sheng-ze HU, and GUO Cheng. 2018. Chinese news hot subtopic discovery and recommendation method based on key phrase and the lda model. *DEStech Transactions on Engineering and Technology Research*, (ecar).
- Steven Bird. 2006. Nltk: the natural language toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.
- Gerlof Bouma. 2009. Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, pages 31–40.
- Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-graber, and David Blei. 2009. [Reading tea leaves: How humans interpret topic models](#). In *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc.
- Pattarawat Chormai, Ponrawee Prasertsom, Jin Cheevaprawatdomrong, and Attapol Rutherford. 2020. Syllable-based neural thai word segmentation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4619–4637.
- Amina Chouigui, Oussama Ben Khiroun, and Bilel Elayeb. 2017. Ant corpus: an arabic news text collection for textual classification. In *2017 IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA)*, pages 135–142. IEEE.
- Ahmed El-Kishky, Yanglei Song, Chi Wang, Clare R. Voss, and Jiawei Han. 2014. [Scalable topical phrase mining from text corpora](#). *Proc. VLDB Endow.*, 8(3):305–316.



- Jey Han Lau, Timothy Baldwin, and David Newman. 2013. On collocations and topic models. *ACM Transactions on Speech and Language Processing (TSLP)*, 10(3):1–14.
- Jey Han Lau, David Newman, and Timothy Baldwin. 2014. Machine reading tea leaves: Automatically evaluating topic coherence and topic model quality. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 530–539.
- Bing Li, Xiaochun Yang, Rui Zhou, Bin Wang, Chengfei Liu, and Yanchun Zhang. 2018. An efficient method for high quality and cohesive topical phrase mining. *IEEE Transactions on Knowledge and Data Engineering*, 31(1):120–137.
- Robert Lindsey, William Headden, and Michael Stipicevic. 2012. A phrase-discovering topic model using hierarchical pitman-yor processes. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 214–222.
- Christopher Manning and Hinrich Schütze. 1999. *Foundations of statistical natural language processing*. MIT press.
- Chandler May, Ryan Cotterell, and Benjamin Van Durme. 2016. An analysis of lemmatization on topic models of morphologically rich language. *arXiv preprint arXiv:1608.03995*.
- Andrew Kachites McCallum. 2002. Mallet: A machine learning for language toolkit. [Http://mallet.cs.umass.edu](http://mallet.cs.umass.edu).
- Paul McCann. 2020. [fugashi, a tool for tokenizing Japanese in python](#). In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 44–51, Online. Association for Computational Linguistics.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- David Mimno, Hanna Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. Optimizing semantic coherence in topic models. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 262–272.
- Ossama Obeid, Nasser Zalmout, Salam Khalifa, Dima Taji, Mai Oudah, Bashar Alhafni, Go Inoue, Fadhil Eryani, Alexander Erdmann, and Nizar Habash. 2020. [CAMEL tools: An open source python toolkit for Arabic natural language processing](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 7022–7032, Marseille, France. European Language Resources Association.
- Eunjeong L. Park and Sungzoon Cho. 2014. Konlpy: Korean natural language processing in python. In *Proceedings of the 26th Annual Conference on Human & Cognitive Language Technology*, Chuncheon, Korea.
- Thomas Proisl and Peter Uhrig. 2016. [SoMaJo: State-of-the-art tokenization for German web and social media texts](#). In *Proceedings of the 10th Web as Corpus Workshop (WAC-X) and the EmpiriST Shared Task*, pages 57–62, Berlin. Association for Computational Linguistics (ACL).
- Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA.
- Peter J Rousseeuw. 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65.
- Evan Sandhaus. 2008. The new york times annotated corpus. *Linguistic Data Consortium, Philadelphia*, 6(12):e26752.
- Alexandra Schofield and David Mimno. 2016. [Comparing apples to apple: The effects of stemmers on topic models](#). *Transactions of the Association for Computational Linguistics*, 4:287–300.
- Huihsin Tseng, Pi-Chuan Chang, Galen Andrew, Dan Jurafsky, and Christopher D Manning. 2005. A conditional random field word segmenter for sghan bakeoff 2005. In *Proceedings of the fourth SIGHAN workshop on Chinese language Processing*.
- Hanna M. Wallach, Iain Murray, Ruslan Salakhutdinov, and David Mimno. 2009. [Evaluation methods for topic models](#). In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, page 1105–1112, New York, NY, USA. Association for Computing Machinery.
- Minmei Wang, Bo Zhao, and Yihua Huang. 2016. Ptr: phrase-based topical ranking for automatic keyphrase extraction in scientific publications. In *International Conference on Neural Information Processing*, pages 120–128. Springer.
- Zhiguo Yu, Todd R Johnson, and Ramakanth Kavuluru. 2013. Phrase based topic modeling for semantic information processing in biomedicine. In *2013 12th International Conference on Machine Learning and Applications*, volume 1, pages 440–445. IEEE.