

# Topical N-grams: Phrase and Topic Discovery, with an Application to Information Retrieval

Xuerui Wang, Andrew McCallum, Xing Wei  
University of Massachusetts  
140 Governors Dr, Amherst, MA 01003  
{xuerui, mccallum, xwei}@cs.umass.edu

## Abstract

*Most topic models, such as latent Dirichlet allocation, rely on the bag-of-words assumption. However, word order and phrases are often critical to capturing the meaning of text in many text mining tasks. This paper presents topical  $n$ -grams, a topic model that discovers topics as well as topical phrases. The probabilistic model generates words in their textual order by, for each word, first sampling a topic, then sampling its status as a unigram or bigram, and then sampling the word from a topic-specific unigram or bigram distribution. Thus our model can model “white house” as a special meaning phrase in the ‘politics’ topic, but not in the ‘real estate’ topic. Successive bigrams form longer phrases. We present experiments showing meaningful phrases and more interpretable topics from the NIPS data and improved information retrieval performance on a TREC collection.*

## 1 Introduction

Although the bag-of-words assumption is prevalent in document classification and topic models, the great majority of natural language processing methods represent word order, including  $n$ -gram language models for speech recognition, finite-state models for information extraction and context-free grammars for parsing. Word order is not only important for syntax, but also important for lexical meaning. A collocation is a phrase with meaning beyond the individual words. For example, the phrase “white house” carries a special meaning beyond the appearance of its individual words, whereas “yellow house” does not. Note, however, that whether or not a phrase is a collocation may depend on the topic context. In the context of a document about real estate, “white house” may not be a collocation.

Most topic models such as latent Dirichlet allocation (LDA) [1], however, assume that words are generated independently from each other, i.e., under the bag-of-words assumption. Adding phrases increases the model’s complexity, but it could be useful in certain contexts.

Assume that we conduct topic analysis on a large collection of research papers. The acknowledgment sections of

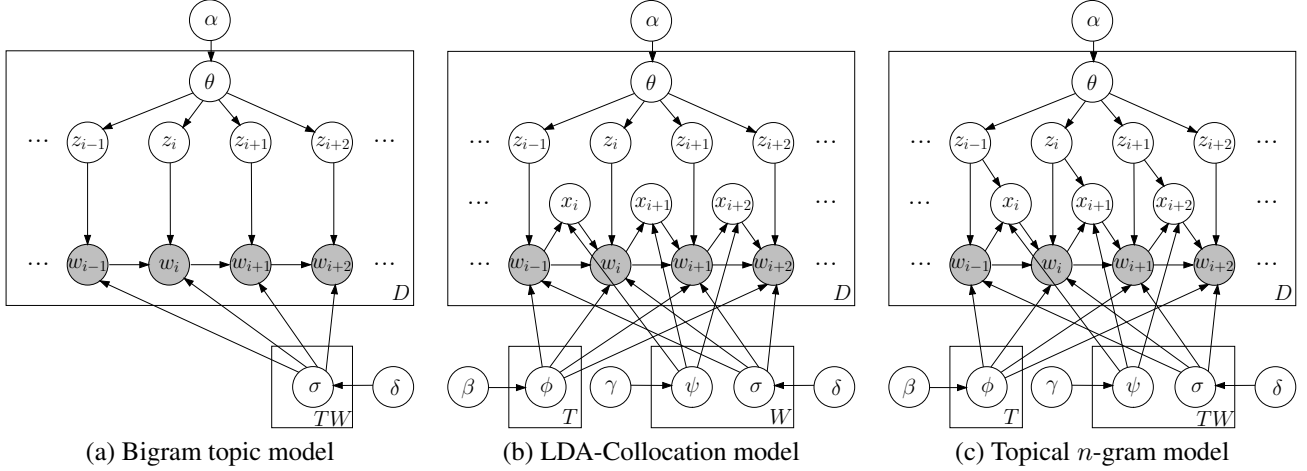
research papers have a distinctive vocabulary. Not surprisingly, we would end up with a particular topic on acknowledgment, since the papers’ acknowledgment sections are not tightly coupled with the content of papers. One might therefore expect to find words such as “thank”, “support” and “grant” in a single topic. One might be very confused, however, to find words like “health” and “science” in the same topic, unless they are presented in context: “National Institutes of Health” and “National Science Foundation”.

Phrases often have specialized meaning, but not always. For example, “neural networks” is considered a phrase because of its frequent use as a fixed expression. However, it specifies two distinct concepts: biological neural networks in neuroscience and artificial neural networks in modern usage. Without consulting the context in which it is located, it is hard to determine its actual meaning. In many situations, topic is very useful to accurately capture the meaning.

In this paper, we propose a new topical  $n$ -gram (TNG) model that automatically determines unigram words and phrases based on context and assign mixture of topics to both individual words and  $n$ -gram phrases. The ability to form phrases only where appropriate is unique to our model, distinguishing it from the traditional collocation discovery methods discussed in Section 3, where a *discovered* phrase is always treated as a *collocation* regardless of the context (which would possibly make us incorrectly conclude that “white house” remains a collocation in a document about real estate). Thus, TNG is not only a topic model that uses phrases, but also help linguists discover meaningful phrases in right context, in a completely probabilistic manner. We show examples of extracted phrases and more interpretable topics on the NIPS data, and in a text mining application, we present better information retrieval performance on an ad-hoc retrieval task over a TREC collection.

## 2 $N$ -gram based Topic Models

Before presenting our topical  $n$ -gram model, we first describe two related  $n$ -gram models. Notation used in this paper is listed in Table 1, and the graphical models are showed in Figure 1. For simplicity, all the models discussed in this section make the 1<sup>st</sup> order Markov assumption, that is, they



**Figure 1. Three  $n$ -gram models ( $D$ : # of documents;  $T$ : # of topics;  $W$ : # of unique words)**

SYMBOL	DESCRIPTION
$T$	number of topics
$D$	number of documents
$W$	number of unique words
$N_d$	number of word tokens in document $d$
$z_i^{(d)}$	the topic associated with the $i^{th}$ token in the document $d$
$x_i^{(d)}$	the bigram status between the $(i-1)^{th}$ token and $i^{th}$ token in the document $d$
$w_i^{(d)}$	the $i^{th}$ token in document $d$

**Table 1. Notation used in this paper**

are actually bigram models. However, all the models have the ability to “model” higher order  $n$ -grams ( $n > 2$ ) by concatenating consecutive bigrams.

### 2.1 Bigram Topic Model (BTM)

Recently, Wallach develops a bigram topic model [14] on the basis of the hierarchical Dirichlet language model [10], by incorporating the concept of topic into bigram models. We assume a dummy word  $w_0$  existing at the beginning of each document. The graphical model presentation of this model is shown in Figure 1(a).

### 2.2 LDA Collocation Model (LDACOL)

Starting from the LDA topic model, the LDA collocation model [13] introduces a new set of random variables (for bigram status)  $x$  ( $x_i = 1$ :  $w_{i-1}$  and  $w_i$  form a bigram;  $x_i = 0$ : they do not) that denote if a bigram can be formed with the previous token, in addition to the two sets of random variables  $z$  and  $w$  in LDA. Thus, it has the power to decide if to generate a bigram or a unigram. At this aspect, it is more realistic than the bigram topic model which always generates bigrams. We assume the status variable  $x_1$  is observed, and only a unigram is allowed at the beginning

of a document. If we want to put more constraints into the model (e.g., no bigram is allowed for sentence/paragraph boundary; only a unigram can be considered for the next word after a stop word is removed; etc.), we can assume that the corresponding status variables are observed as well. Its graphical model presentation is shown in Figure 1(b).

### 2.3 Topical $N$ -gram Model (TNG)

The topical  $n$ -gram model (TNG) is not a pure addition of the bigram topic model and LDA collocation model. One of the key contributions of our model is to make it possible to decide whether to form a bigram for the same two consecutive word tokens depending on their nearby context (i.e., co-occurrences). As in the LDA collocation model, we may assume some  $x$ ’s are observed for the same reason as we discussed in Section 2.2. The graphical model presentation of this model is shown in Figure 1(c). Note that our model is a more powerful generalization of BTM and of LDACOL. Both BTM (by setting all  $x$ ’s to 1) and LDACOL (by making  $\sigma$  conditioned on previous word only) are the special cases of our TNG model.

Its generative process can be described as follows: 1) draw Discrete  $\phi_z$  from Dirichlet  $\beta$  for each topic  $z$ ; 2) draw Bernoulli  $\psi_{zw}$  from Beta  $\gamma$  for each topic  $z$  and each word  $w$ ; 3) draw Discrete  $\sigma_{zw}$  from Dirichlet  $\delta$  for each topic  $z$  and each word  $w$ ; 4) for each document  $d$ , draw Discrete  $\theta^{(d)}$  from Dirichlet  $\alpha$ ; then for each word  $w_i^{(d)}$  in document  $d$ : 4a) draw  $x_i^{(d)}$  from Bernoulli  $\psi_{z_{i-1}^{(d)} w_{i-1}^{(d)}}$ ; 4b) draw  $z_i^{(d)}$  from Discrete  $\theta^{(d)}$ ; and 4c) draw  $w_i^{(d)}$  from Discrete  $\sigma_{z_i^{(d)} w_{i-1}^{(d)}}$  if  $x_i^{(d)} = 1$ ; else draw  $w_i^{(d)}$  from Discrete  $\phi_{z_i^{(d)}}$ .

As shown in the above, the topic assignments for the two terms in a bigram are not required to be identical. In this paper, we will use the topic of the last term as the topic of the phrase for simplicity, since long noun phrases do truly

sometimes have components indicative of different topics, and its last noun is usually the “head noun”. Alternatively, we could enforce consistency in the model with ease, by simply adding two more sets of arrows ( $z_{i-1} \rightarrow z_i$  and  $x_i \rightarrow z_i$ ). Accordingly, we could substitute Step 4b) in the above generative process with “draw  $z_i^{(d)}$  from Discrete  $\theta^{(d)}$  if  $x_i^{(d)} = 1$ ; else let  $z_i^{(d)} = z_{i-1}^{(d)}$ ,” In this way, a word has the option to inherit a topic assignment from its previous word if they form a bigram phrase. However, from our experimental results, the first choice yields better performance. We will focus on the model shown in Figure 1(c).

In state-of-the-art hierarchical Bayesian models such as latent Dirichlet allocation, exact inference over hidden topic variables is typically intractable due to the large number of latent variables and parameters in the models. We use Gibbs sampling to conduct approximate inference in this paper. To reduce the uncertainty introduced by  $\theta, \phi, \psi$ , and  $\sigma$ , we could integrate them out with no trouble because of the conjugate prior setting in our model. Starting from the joint distribution  $P(\mathbf{w}, \mathbf{z}, \mathbf{x} | \alpha, \beta, \gamma, \delta)$ , we can work out the conditional probabilities  $P(z_i^{(d)}, x_i^{(d)} | \mathbf{z}_{-i}^{(d)}, \mathbf{x}_{-i}^{(d)}, \mathbf{w}, \alpha, \beta, \gamma, \delta)$  conveniently using Bayes rule, where  $\mathbf{z}_{-i}^{(d)}$  denotes the topic assignments for all word tokens except word  $w_i^{(d)}$ , and  $\mathbf{x}_{-i}^{(d)}$  represents the bigram status for all tokens except word  $w_i^{(d)}$ . During Gibbs sampling, we draw the topic assignment  $z_i^{(d)}$  and the bigram status  $x_i^{(d)}$  iteratively for each word token  $w_i^{(d)}$  according to the conditional probability distribution:

$$P(z_i^{(d)}, x_i^{(d)} | \mathbf{z}_{-i}^{(d)}, \mathbf{x}_{-i}^{(d)}, \mathbf{w}, \alpha, \beta, \gamma, \delta) \propto (\gamma_{x_i^{(d)}} + p_{z_{i-1}^{(d)} w_{i-1}^{(d)} x_i} - 1)(\alpha_{z_i^{(d)}} + q_{dz_i^{(d)}} - 1) \times \begin{cases} \frac{\beta_{w_i^{(d)}} + n_{z_i^{(d)} w_i^{(d)}} - 1}{\sum_{v=1}^W (\beta_v + n_{z_i^{(d)} v}) - 1} & \text{if } x_i^{(d)} = 0 \\ \frac{\delta_{w_i^{(d)}} + m_{z_i^{(d)} w_{i-1}^{(d)} w_i^{(d)}} - 1}{\sum_{v=1}^W (\delta_v + m_{z_i^{(d)} w_{i-1}^{(d)} v}) - 1} & \text{if } x_i^{(d)} = 1 \end{cases}$$

where  $n_{zw}$  represents how many times word  $w$  is assigned into topic  $z$  as a unigram,  $m_{zdwv}$  represents how many times word  $v$  is assigned to topic  $z$  as the  $2^{nd}$  term of a bigram given the previous word  $w$ ,  $p_{zwk}$  denotes how many times the status variable  $x = k$  (0 or 1) given the previous word  $w$  and the previous word’s topic  $z$ , and  $q_{dz}$  represents how many times a word is assigned to topic  $z$  in document  $d$ . Note all counts here do include the assignment of the token being visited. Simple manipulations give us the posterior estimates of  $\theta, \phi, \psi$ , and  $\sigma$  as follows:

$$\begin{aligned} \hat{\theta}_z^{(d)} &= \frac{\alpha_z + q_{dz}}{\sum_{t=1}^T (\alpha_t + q_{dt})} & \hat{\phi}_{zw} &= \frac{\beta_w + n_{zw}}{\sum_{v=1}^W (\beta_v + n_{zv})} \\ \hat{\psi}_{zwk} &= \frac{\gamma_k + p_{zwk}}{\sum_{k=0}^1 (\gamma_k + p_{zwk})} & \hat{\sigma}_{zwv} &= \frac{\delta_v + m_{zwv}}{\sum_{v=1}^W (\delta_v + m_{zwv})} \end{aligned} \quad (1)$$

### 3 Related Work

Collocation has long been studied by lexicographers and linguists in various ways. Traditional collocation discov-

ery methods range from frequency to variance, to hypothesis testing, to mutual information. The simplest method is counting. A small amount of linguistic knowledge (a part-of-speech filter) has been combined with frequency [9] to discover surprisingly meaningful phrases. Variance based collocation discovery [12] considers collocations in a more flexible way than fixed phrases. However, high frequency and low variance can be accidental. Hypothesis testing can be used to assess whether or not two words occur together more often than chance. Many statistical tests have been explored, for example,  $t$ -test [3],  $\chi^2$  test [2], and likelihood ratio test [4]. More recently, an information-theoretically motivated method for collocation discovery is utilizing mutual information [8].

The hierarchical Dirichlet language model [10] is closely related to the bigram topic model [14]. The probabilistic view of smoothing in language models shows how to take advantage of a bigram model in a Bayesian way.

The main stream of topic modeling has gradually gained a probabilistic flavor as well in the past decade. One of the most popular topic model, latent Dirichlet allocation (LDA), which makes the bag-of-words assumption, has made a big impact in the fields of natural language processing, statistical machine learning and text mining. Three models we discussed in Section 2 all contain an LDA component that is responsible for the topic part.

In our point of view, the HMMLDA model [7] is the first attack to word dependency in the topic modeling framework. The authors present HMMLDA as a generative composite model that takes care of both short-range syntactic dependencies and long-range semantic dependencies between words; its syntactic part is a hidden Markov model and the semantic component is a topic model (LDA). Interesting results based on this model are shown on tasks such as part-of-speech tagging and document classification.

### 4 Experimental Results

We apply the topical  $n$ -gram model to the NIPS dataset that consists of the 13 years of proceedings from 1987 to 1999 Neural Information Processing Systems (NIPS) Conferences. The dataset contains 1,740 research papers, 13,649 unique words, and 2,301,375 word tokens in total. Topics found from a 50-topic run on the NIPS dataset are shown in Table 2 as anecdotal evidence, with comparison to the corresponding closest LDA topics.

The “Reinforcement Learning” topic provides an extremely salient summary of the corresponding research area. The LDA topic assembles many common words used in reinforcement learning, but in its word list, there are quite a few generic words (such as “function”, “dynamic”, “decision”) that are common and highly probable in many other topics as well. In TNG, we can find that these generic words are associated with other words to form  $n$ -gram phrases

Reinforcement Learning			Human Receptive System		
LDA	$n$ -gram (2+)	$n$ -gram (1)	LDA	$n$ -gram (2+)	$n$ -gram (1)
state	reinforcement learning	action	motion	receptive field	motion
learning	optimal policy	policy	visual	spatial frequency	spatial
policy	dynamic programming	reinforcement	field	temporal frequency	visual
action	optimal control	states	position	visual motion	receptive
reinforcement	function approximator	actions	figure	motion energy	response
states	prioritized sweeping	function	direction	tuning curves	direction
time	finite-state controller	optimal	fields	horizontal cells	cells
optimal	learning system	learning	eye	motion detection	figure
actions	reinforcement learning rl	reward	location	preferred direction	stimulus
function	function approximators	control	retina	visual processing	velocity
algorithm	markov decision problems	agent	receptive	area mt	contrast
reward	markov decision processes	q-learning	velocity	visual cortex	tuning
step	local search	goal	vision	light intensity	moving
dynamic	state-action pair	space	moving	directional selectivity	model
control	markov decision process	step	system	high contrast	temporal
sutton	belief states	environment	flow	motion detectors	responses
rl	stochastic policy	system	edge	spatial phase	orientation
decision	action selection	problem	center	moving stimuli	light
algorithms	upright position	steps	light	decision strategy	stimuli
agent	reinforcement learning methods	transition	local	visual stimuli	cell

Speech Recognition			Support Vector Machines		
LDA	$n$ -gram (2+)	$n$ -gram (1)	LDA	$n$ -gram (2+)	$n$ -gram (1)
recognition	speech recognition	speech	kernel	support vectors	kernel
system	training data	word	linear	test error	training
word	neural network	training	vector	support vector machines	support
face	error rates	system	support	training error	margin
context	neural net	recognition	set	feature space	svm
character	hidden markov model	hmm	nonlinear	training examples	solution
hmm	feature vectors	speaker	data	decision function	kernels
based	continuous speech	performance	algorithm	cost functions	regularization
frame	training procedure	phoneme	space	test inputs	adaboost
segmentation	continuous speech recognition	acoustic	pca	kkt conditions	test
training	gamma filter	words	function	leave-one-out procedure	data
characters	hidden control	context	problem	soft margin	generalization
set	speech production	systems	margin	bayesian transduction	examples
probabilities	neural nets	frame	vectors	training patterns	cost
features	input representation	trained	solution	training points	convex
faces	output layers	sequence	training	maximum margin	algorithm
words	training algorithm	phonetic	svm	strictly convex	working
frames	test set	speakers	kernels	regularization operators	feature
database	speech frames	mlp	matrix	base classifiers	sv
mlp	speaker dependent	hybrid	machines	convex optimization	functions

**Table 2. The four topics from a 50-topic run of TNG on the NIPS data with their closest counterparts from LDA. The Title above the word list of each topic is our own summary of the topic. To better illustrate the difference between TNG and LDA, we list the  $n$ -grams ( $n > 1$ ) and unigrams separately for TNG. Each topic is shown with the 20 sorted highest-probability words. TNG produces clearer word list for each topic by associating many generic words with other words to form  $n$ -gram phrases.**

(such as “markov decision process”, etc.) that are only highly probable in reinforcement learning. More importantly, by forming  $n$ -gram phrases, the unigram word list produced by TNG is also cleaner. For example, because of the prevalence of generic words in LDA, highly related words (such as “q-learning” and “goal”) are not ranked high enough to be shown in the top 20 word list. On the contrary, they are ranked very high in the TNG’s unigram word list.

In the other three topics (Table 2), we can find similar phenomena as well. For example, in “Human Receptive System”, some generic words (such as “field”, “receptive”) are actually the components of the popular phrases in this area as shown in the TNG model. “system” is ranked high in LDA, but almost meaningless, and on the other hand, it does not appear in the top word lists of TNG. Some extremely related words (such as “spatial”), ranked very high in TNG, are absent in LDA’s top word list. In “Speech Recognition”,

the dominating generic words (such as “context”, “based”, “set”, “probabilities”, “database”) make the LDA topic less understandable than even just TNG’s unigram word list.

In many situations, a crucially related word might be not mentioned enough to be clearly captured in LDA, on the other hand, it would become very salient as a phrase due to the relatively stronger co-occurrence pattern in an extremely sparse setting for phrases. The “Support Vector Machines” topic provides such an example. We can imagine that “kkt” will be mentioned no more than a few times in a typical NIPS paper, and it probably appears only as a part of the phrase “kkt conditions”. TNG satisfyingly captures it successfully as a highly probable phrase in the SVM topic.

As discussed before, higher-order  $n$ -grams ( $n > 2$ ) can be approximately modeled by concatenating consecutive bi-grams in the TNG model, as shown in Table 2 (such as “hidden markov model” and “support vector machines”, etc.).

## 4.1 Ad-hoc Retrieval

Traditional information retrieval (IR) models usually represent text with bags-of-words assuming that words occur independently, which is not exactly appropriate to natural language. To address this problem, researchers have been working on capturing word dependencies. There are mainly two types of dependencies being studied and shown to be effective: 1) topical (semantic) dependency, which is also called long-distance dependency. Two words are considered dependent when their meanings are related and they co-occur often, such as “fruit” and “apple”. Among models capturing semantic dependency, the LDA-based document models [15] are state-of-the-art. For IR applications, a major advantage of applying topic models to document expansion, compared to online query expansion in pseudo relevance feedback, is that they can be trained offline, thus more efficient in handling a new query; 2) phrase dependency, also called short-distance dependency. As reported in literature, retrieval performance can be boosted if the similarity between a user query and a document is calculated by common phrases instead of common words [6, 5, 11]. Most research on phrases in information retrieval has employed an independent collocation discovery module. In this way, a phrase can be indexed exactly as an ordinary word.

The topical  $n$ -gram model automatically and simultaneously takes cares of both semantic co-occurrences and phrases. Also, it does not need a separate module for phrase discovery, and everything can be seamlessly integrated into the language modeling framework, which is one of the most popular statistically principled approaches to IR.

The SJMN dataset, taken from TREC with standard queries 51-150 that are taken from the *title* field of TREC topics, covers materials from San Jose Mercury News in 1991. In total, the SJMN dataset we use contains 90,257 documents, 150,714 unique words, and 21,156,378 tokens, which is order of magnitude larger than the NIPS dataset. In TNG, a word distribution for each document can be calculated, which thus can be viewed as a document model, and the likelihood of generating a query can be computed to rank documents, which is the basic idea in the query likelihood (QL) model in IR. In the query likelihood model, each document is scored by the likelihood of its model generating a query  $Q$ ,  $P_{LM}(Q|d)$ . Let the query  $Q = (q_1, q_2, \dots, q_{L_Q})$ . Under the bag-of-words assumption,  $P_{LM}(Q|d) = \prod_{i=1}^{L_Q} P(q_i|d)$ , which is often specified by the document model with Dirichlet smoothing,  $P_{LM}(q|d) = \frac{N_d}{N_d + \mu} P_{ML}(q|d) + (1 - \frac{N_d}{N_d + \mu}) P_{ML}(q|coll)$ , where  $N_d$  is the length of document  $d$ ,  $P_{ML}(q|d)$  and  $P_{ML}(q|coll)$  are the maximum likelihood (ML) estimates of a query term  $q$  generated in document  $d$  and in the entire collection, respectively, and  $\mu = 1000$  is the Dirichlet smoothing prior.

To calculate the QL from the TNG model within the language modeling framework, we need to sum over the

topic variable and bigram status variable for each token. Given the posterior estimates  $\hat{\theta}$ ,  $\hat{\phi}$ ,  $\hat{\psi}$ , and  $\hat{\sigma}$  (Eqn. 1), the QL of query  $Q$  given document  $d$ ,  $P_{TNG}(Q|d)$  can be calculated<sup>1</sup> as  $P_{TNG}(Q|d) = \prod_{i=1}^{L_Q} P_{TNG}(q_i|q_{i-1}, d)$ , where  $P_{TNG}(q_i|q_{i-1}, d) = \sum_{z_i=1}^T (P(x_i = 0|\hat{\psi}_{q_{i-1}})P(q_i|\hat{\phi}_{z_i}) + P(x_i = 1|\hat{\psi}_{q_{i-1}})P(q_i|\hat{\sigma}_{z_i q_{i-1}}))P(z_i|\hat{\theta}^{(d)})$ , and,  $P(x_i|\hat{\psi}_{q_{i-1}}) = \sum_{z_{i-1}=1}^T P(x_i|\hat{\psi}_{z_{i-1} q_{i-1}})P(z_{i-1}|\hat{\theta}^{(d)})$ . Note in the calculation, the bag-of-words assumption is not made any more.

Similar to the method in [15], we can combine the query likelihood from the basic language model and the likelihood from the TNG model in various ways. One can combine them at query level, i.e.,  $P(Q|d) = \lambda P_{LM}(Q|d) + (1 - \lambda)P_{TNG}(Q|d)$ , where  $\lambda$  is a weighting factor between the two likelihoods.

Alternatively, under first order Markov assumption,  $P(Q|d) = P(q_1|d) \prod_{i=2}^{L_Q} P(q_i|q_{i-1}, d)$ , and one can combine the query likelihood at query term level (used in this paper), that is,  $P(q_i|q_{i-1}, d) = \lambda P_{LM}(q_i|d) + (1 - \lambda)P_{TNG}(q_i|q_{i-1}, d)$ . The query likelihood  $P(Q|d)$  for the BTM and LDACOL models can be calculated similarly and more simply, considering both models are special cases of our TNG model.

## 4.2 Comparison of BTM, LDACOL and TNG on TREC Ad-hoc Retrieval

In this section, we compare the IR performance of the three  $n$ -gram based topic models on the SJMN dataset, as shown in Table 3. For a fair comparison, the weighting factor  $\lambda$  (reported in Table 3) are independently chosen to get the best performance from each model. Under the Wilcoxon test with 95% confidence, TNG significantly outperforms BTM and LDACOL on this standard retrieval task.

Space limitations prevent us from presenting the results for all queries, but it is interesting to see that different models are good at quite different queries. For some queries (such as No. 117 and No. 138), TNG and BTM perform similarly, and better than LDACOL, and for some other queries (such as No. 110 and No. 150), TNG and LDACOL perform similarly, and better than BTM. There are also queries (such as No. 061 and No. 130) for which TNG performs better than both BTM and LDACOL. We believe that they are clear empirical evidence that our TNG model are more generic and powerful than BTM and LDACOL.

It is true that for certain queries (such as No. 069 and No. 146), TNG performs worse than BTM and LDACOL, but we notice that all models perform badly on these queries and the behaviors are more possibly due to randomness.

## 5 Conclusions

In this paper, we have presented the topical  $n$ -gram model. The TNG model automatically determines to form

<sup>1</sup>A dummy  $q_0$  is assumed at the beginning of every query, for the convenience of mathematical presentation.

No.	Query	TNG	BTM	Change	LDACOL	Change
061	Israeli Role in Iran-Contra Affair	0.1635	0.1104	-32.47%	0.1316	-19.49%
069	Attempts to Revive the SALT II Treaty	0.0026	0.0071	172.34%	0.0058	124.56%
110	Black Resistance Against the South African Government	0.4940	0.3948	-20.08%	0.4883	-1.16%
117	Capacity of the U.S. Cellular Telephone Network	0.2801	0.3059	9.21%	0.1999	-28.65%
130	Jewish Emigration and U.S.-USSR Relations	0.2087	0.1746	-16.33%	0.1765	-15.45%
138	Iranian Support for Lebanese Hostage-takers	0.4398	0.4429	0.69%	0.3528	-19.80%
146	Negotiating an End to the Nicaraguan Civil War	0.0346	0.0682	97.41%	0.0866	150.43%
150	U.S. Political Campaign Financing	0.2672	0.2323	-13.08%	0.2688	0.59%
	All Queries	0.2122	0.1996	-5.94%*	0.2107	-0.73%*

**Table 3. Comparison of the bigram topic model ( $\lambda = 0.7$ ), LDA collocation model ( $\lambda = 0.9$ ) and the topical  $n$ -gram Model ( $\lambda = 0.8$ ) on TREC retrieval performance (average precision). \* indicates statistically significant differences in performance with 95% confidence according to the Wilcoxon test. TNG performs significantly better than other two models overall.**

an  $n$ -gram (and further assign a topic) or not, based on its surrounding context. Examples of topics found by TNG are more interpretable than its LDA counterpart. We also demonstrate how TNG can help improve retrieval performance in standard ad-hoc retrieval tasks on TREC collections over its two special-case  $n$ -gram based topic models.

Unlike some traditional phrase discovery methods, the TNG model provides a systematic way to model (topical) phrases and can be seamlessly integrated with many probabilistic frameworks for various tasks such as phrase discovery, ad-hoc retrieval, machine translation, speech recognition and statistical parsing.

To the best of our knowledge, our paper presents the very first application of all three  $n$ -gram based topic models on Gigabyte collections, and a novel way to integrate  $n$ -gram based topic models into the language modeling framework for information retrieval tasks.

## Acknowledgments

This work was supported in part by the Center for Intelligent Information Retrieval and in part by the Defense Advanced Research Projects Agency, through the Department of the Interior, NBC, Acquisition Services Division, under contract #NBCHD030010, and under contract #HR0011-06-C-0023. Any opinions, findings, conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect those of the sponsor.

## References

- [1] D. M. Blei, A. Y. Ng, and M. J. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [2] K. Church and W. Gale. Concordances for parallel text. In *Proceedings of the Seventh Annual Conference of the UW Centre for the New OED and Text Research*, pages 40–62, 1991.
- [3] K. Church and P. Hanks. Word association norms, mutual information and lexicography. In *Proceedings of the 27th*

*Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 76–83, 1989.

- [4] T. E. Dunning. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61–74, 1993.
- [5] D. A. Evans, K. Ginther-Webster, M. Hart, R. G. Lefferts, and I. A. Monarch. Automatic indexing using selective NLP and first-order thesauri. In *Proceedings of Intelligent Multimedia Information Retrieval Systems and Management (RIA0'91)*, pages 624–643, 1991.
- [6] J. Fagan. The effectiveness of a nonsyntactic approach to automatic phrase indexing for document retrieval. *Journal of the American Society for Information Science*, 40(2):115–139, 1989.
- [7] T. Griffiths, M. Steyvers, D. Blei, and J. Tenenbaum. Integrating topics and syntax. In *Advances in Neural Information Processing Systems 17*, 2005.
- [8] J. Hodges, S. Yie, R. Reighart, and L. Boggess. An automated system that assists in the generation of document indexes. *Natural Language Engineering*, 2(2):137–160, 1996.
- [9] J. S. Justeson and S. M. Katz. Technical terminology: some linguistic properties and an algorithm for identification in text. *Natural Language Engineering*, 1:9–27, 1995.
- [10] D. J. C. MacKay and L. Peto. A hierarchical Dirichlet language model. *Natural Language Engineering*, 1(3):1–19, 1994.
- [11] M. Mitra, C. Buckley, A. Singhal, and C. Cardie. An analysis of statistical and syntactic phrases. In *Proceedings of RIAO-97, 5th International Conference*, pages 200–214, Montreal, CA, 1997.
- [12] F. Smadja. Retrieving collocations from text: Xtract. *Computational Linguistics*, 19:143–177, 1993.
- [13] M. Steyvers and T. Griffiths. Matlab topic modeling toolbox 1.3. [http://psiexp.ss.uci.edu/research/programs\\_data/toolbox.htm](http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm), 2005.
- [14] H. Wallach. Topic modeling: beyond bag-of-words. In *Proceedings of the 23rd International Conference on Machine Learning*, 2006.
- [15] X. Wei and W. B. Croft. LDA-based document models for ad-hoc retrieval. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research & Development on Information Retrieval*, 2006.