

Enhancement of Enterprise Search With Neural Network Language Model

Pengchu Zhang, Laritza Saenz and John Mareda
Sandia National Laboratories
November 30, 2015





Problems in Enterprise Search

- ◆ Search engines are basically matching the query terms and the terms in documents, this may result in a large number of false positives;
- ◆ Some terms may be overly represented in search:
 - E.g., "composite materials"
- ◆ Search engines may fail to deliver if the query terms not exist in the enterprise datasets but the datasets have the relevant information, e.g.,
 - Buckyballs vs. fullerene





Current Approaches to Enhance Information finding

- ◆ Boosting some terms;
 - E.g., Composite(+2) materials
- ◆ Synonym to expand query terms;
 - Java vs. coffee
- ◆ Others such as applying various similarity algorithms

In this presentation, we propose to improve information findability with Neural Network Language Models



NNLM in NATURE

From Nature 521, 4346-444 (28 May 2015)

Identific x Search x Yahoo x W Paracon x Expert S x www.na x delivery x

www.nature.com/nature/journal/v521/n7553/pdf/nature14539.pdf

REVIEW

doi:10.1038/nature14539

Deep learning

Yann LeCun^{1,2}, Yoshua Bengio³ & Geoffrey Hinton⁴

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains such as drug discovery and genomics. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep convolutional nets have brought about breakthroughs in processing images, video, speech and audio, whereas recurrent nets have shone light on sequential data such as text and speech.

Machine-learning technology powers many aspects of modern society: from web searches to content filtering on social networks to recommendations on e-commerce websites, and it is increasingly present in consumer products such as cameras and smartphones. Machine-learning systems are used to identify objects in images, transcribe speech into text, match news items, posts or products with users' interests, and select relevant results of search. Increasingly, these applications make use of a class of techniques called deep learning.

Conventional machine-learning techniques were limited in their ability to process natural data in their raw form. For decades, constructing a pattern-recognition or machine-learning system required careful engineering and considerable domain expertise to design a feature extractor that transformed the raw data (such as the pixel values of an image) into a suitable internal representation or feature vector from which the learning subsystem, often a classifier, could detect or classify patterns in the input.

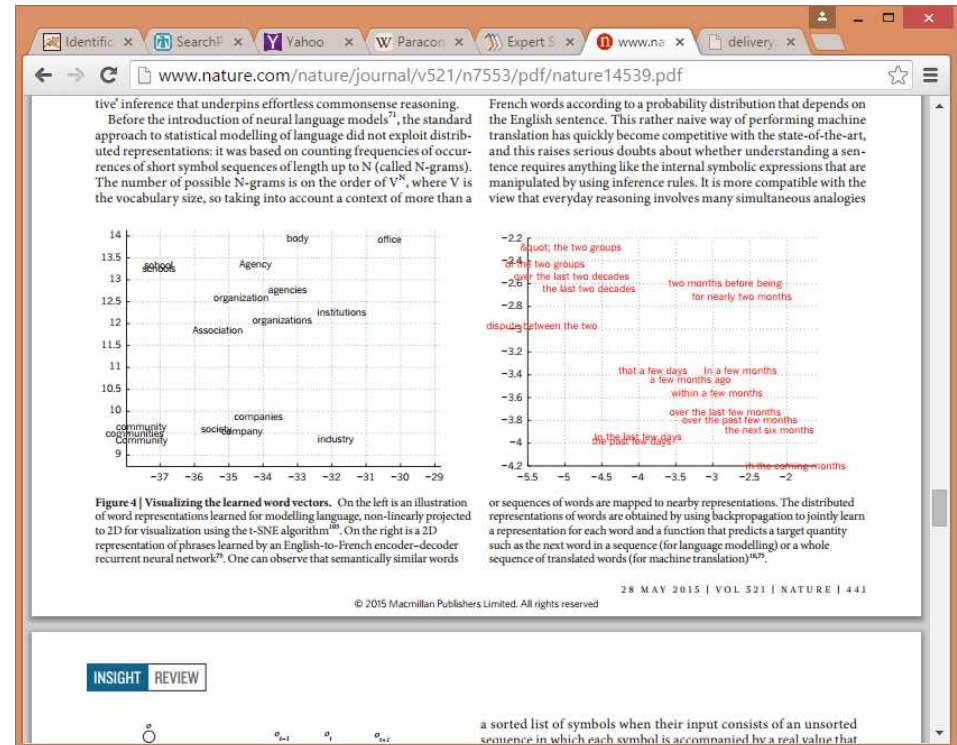
Representation learning is a set of methods that allows a machine to be fed with raw data and to automatically discover the representations needed for detection or classification. Deep-learning methods are representation-learning methods with multiple levels of representation, obtained by composing simple but non-linear modules that each transform the representation at one level (starting with the raw input) into a representation at a higher, slightly more abstract level. With the composition of enough such transformations, very complex functions can be learned. For classification tasks, higher layers of representation amplify aspects of the input that are important for discrimination and suppress irrelevant variations. An image, for example, comes in the form of an array of pixel values, and the learned features in the first layer of representation typically represent the presence or absence of edges at particular orientations and locations in the image. The second layer typically detects motifs by spotting particular arrangements of

intricate structures in high-dimensional data and is therefore applicable to many domains of science, business and government. In addition to beating records in image recognition¹⁻⁴ and speech recognition^{5,6}, it has beaten other machine-learning techniques at predicting the activity of potential drug molecules⁷, analysing particle accelerator data^{8,9}, reconstructing brain circuits¹⁰, and predicting the effects of mutations in non-coding DNA on gene expression and disease^{11,12}. Perhaps more surprisingly, deep learning has produced extremely promising results for various tasks in natural language understanding¹³, particularly topic classification, sentiment analysis, question answering¹⁴ and language translation^{15,16}.

We think that deep learning will have many more successes in the near future because it requires very little engineering by hand, so it can easily take advantage of increases in the amount of available computation and data. New learning algorithms and architectures that are currently being developed for deep neural networks will only accelerate this progress.

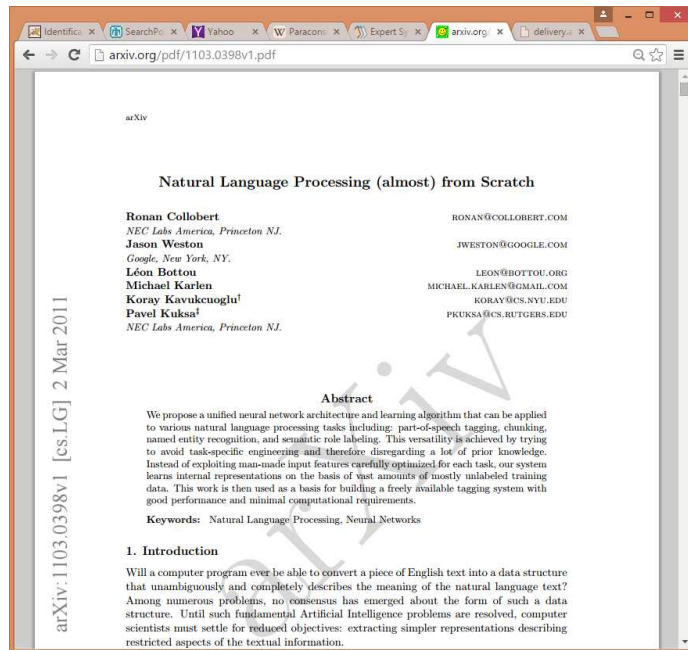
Supervised learning

The most common form of machine learning, deep or not, is supervised learning. Imagine that we want to build a system that can classify images as containing, say, a house, a car, a person or a pet. We first collect a large data set of images of houses, cars, people and pets, each labelled with its category. During training, the machine is shown an image and produces an output in the form of a vector of scores, one for each category. We want the desired category to have the highest score of all categories, but this is unlikely to happen before training. We compute an objective function that measures the error (or distance) between the output scores and the desired pattern of scores. The machine then modifies its internal adjustable parameters to reduce this error. These adjustable parameters, often called weights, are real numbers that can be seen as 'knobs' that define the input-output func-

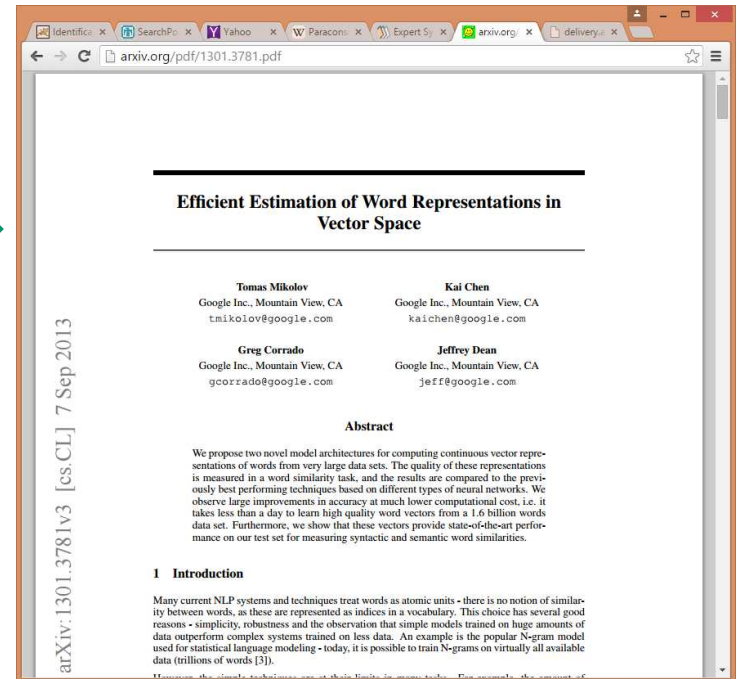


From Theory to Application to Codes

[arXiv:1301.3781v3](https://arxiv.org/abs/1301.3781v3)



[arXiv:1103.0398v1](https://arxiv.org/abs/1103.0398v1)



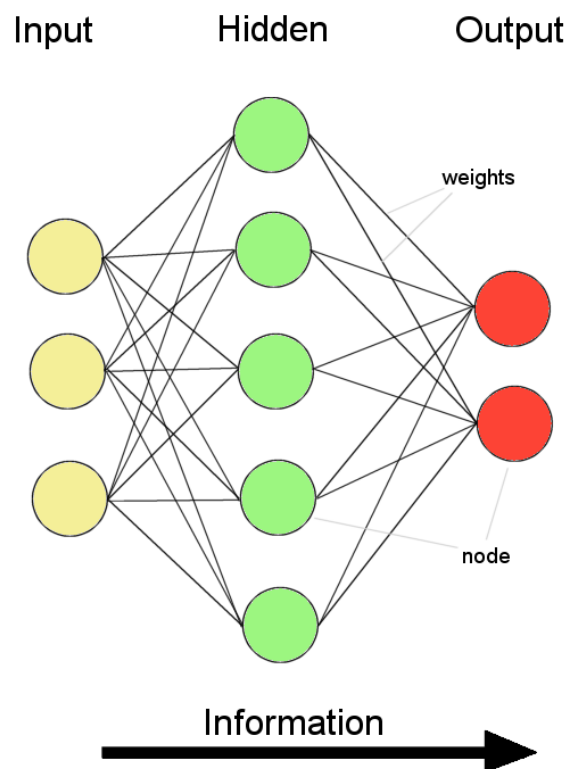
<https://code.google.com/p/word2vec/>

Sandia National Laboratories

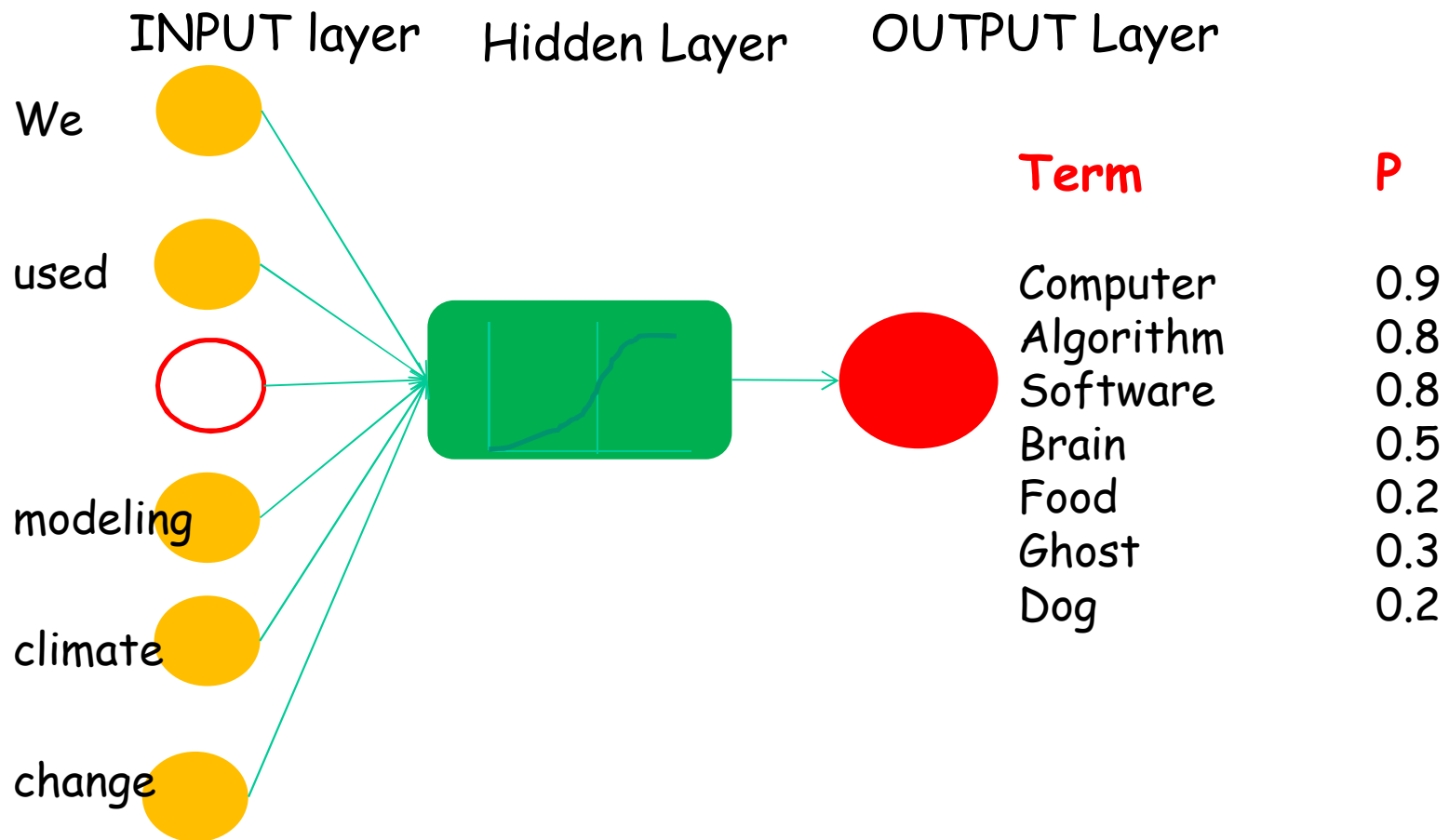


What is Neural Network Language Model?

Neural Network



Concept of Neural Network Language Model



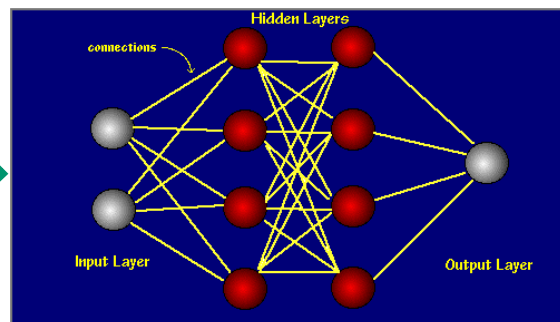
Build and optimize word vectors(Google 2013)

Word in sentence vector

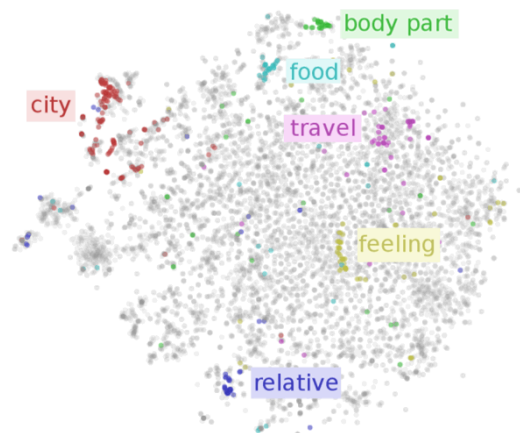
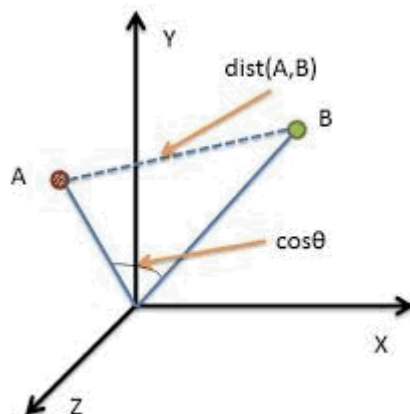
New vector

The
dog
is
walking
in
the
room

(0.12, 0.23, 0.22)
(0.32, 0.27, 0.94)
(0.18, 0.88, 0.45)
(0.23, 0.92, 0.23)
(0.77, 0.25, 0.11)
(0.12, 0.23, 0.22)
(0.41, 0.13, 0.29)



(0.12, 0.23, 0.22)
(0.62, 0.99, 0.14)
(0.18, 0.88, 0.45)
(0.23, 0.92, 0.23)
(0.77, 0.25, 0.11)
(0.12, 0.23, 0.22)
(0.41, 0.13, 0.29)



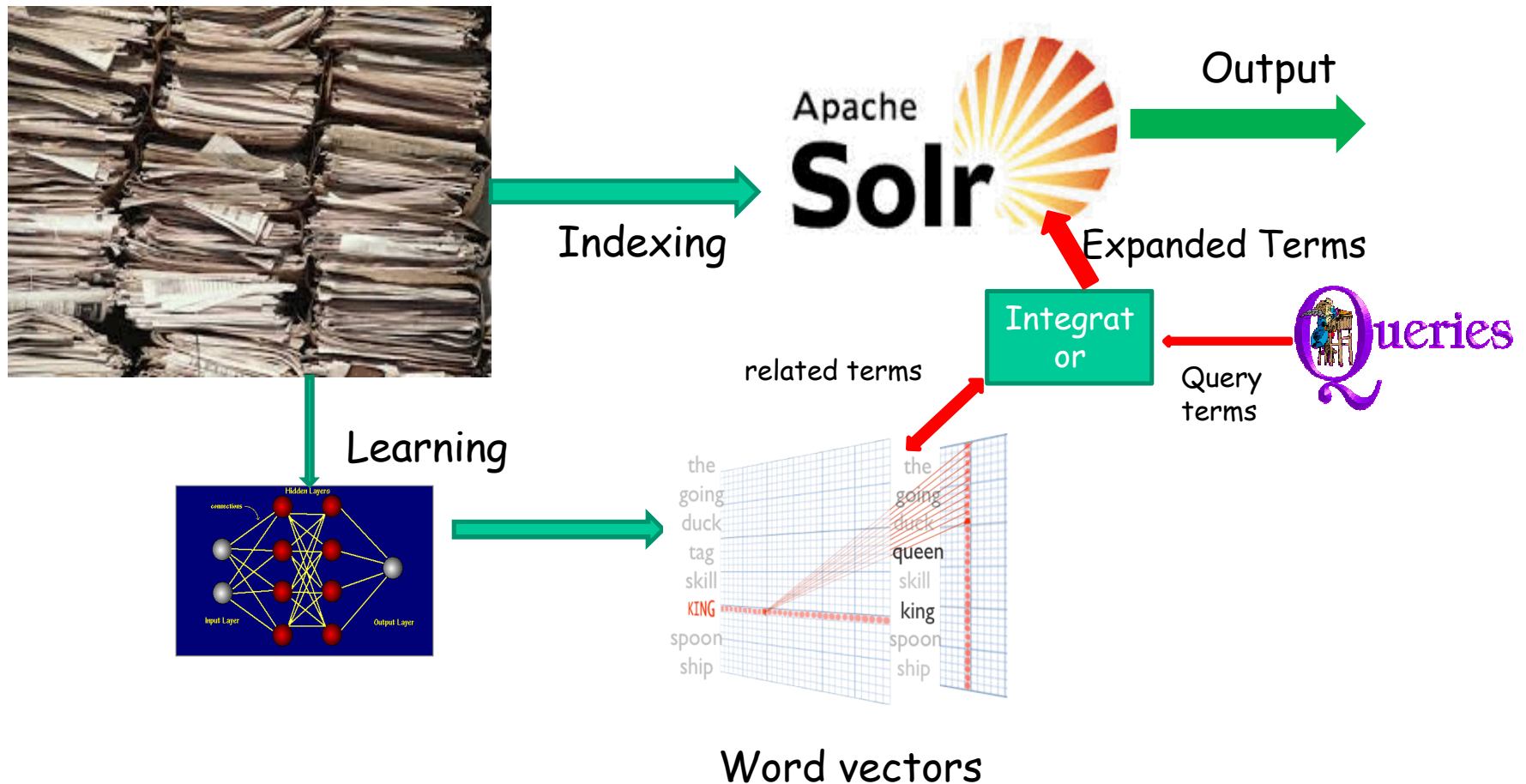


Motivations in Applying Neural Network Language Model (NNLM)

- ◆ NNLM projects the terms into vector space, meaning the words that were used together will stay together at the vector space;
- ◆ NNLM vectorizes the terms that makes it possible to compute the distances between the terms;
- ◆ Most importantly, NNLM learns from natural language semantically and syntactically. With this NNLM meaningfully brings words together without parsing a speech or a document.



Concept for Enhanced Enterprise Search



Note: All images are from online





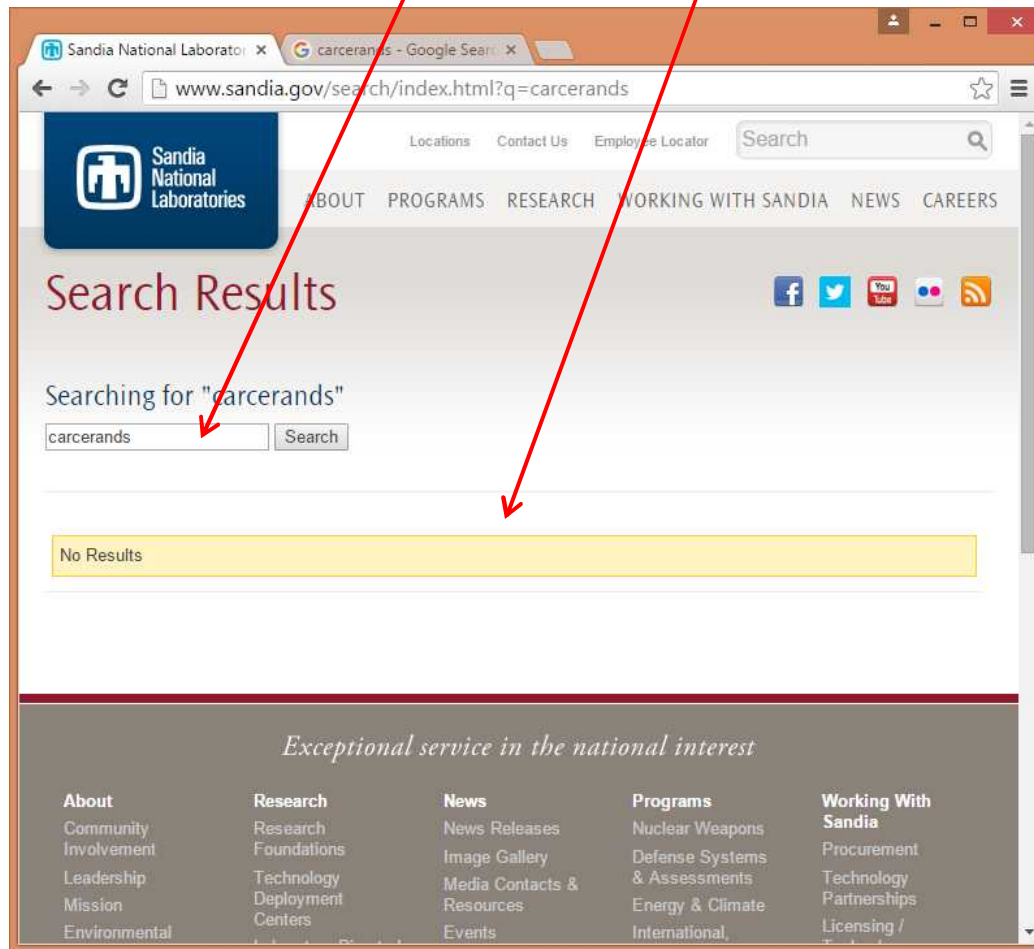
What and How NNLM helps Enterprise Search

- ◆ Deliver the relevant information even the query terms not in enterprise dataset;
- ◆ Increase the quality of information delivered:
 - Relevance
 - Related

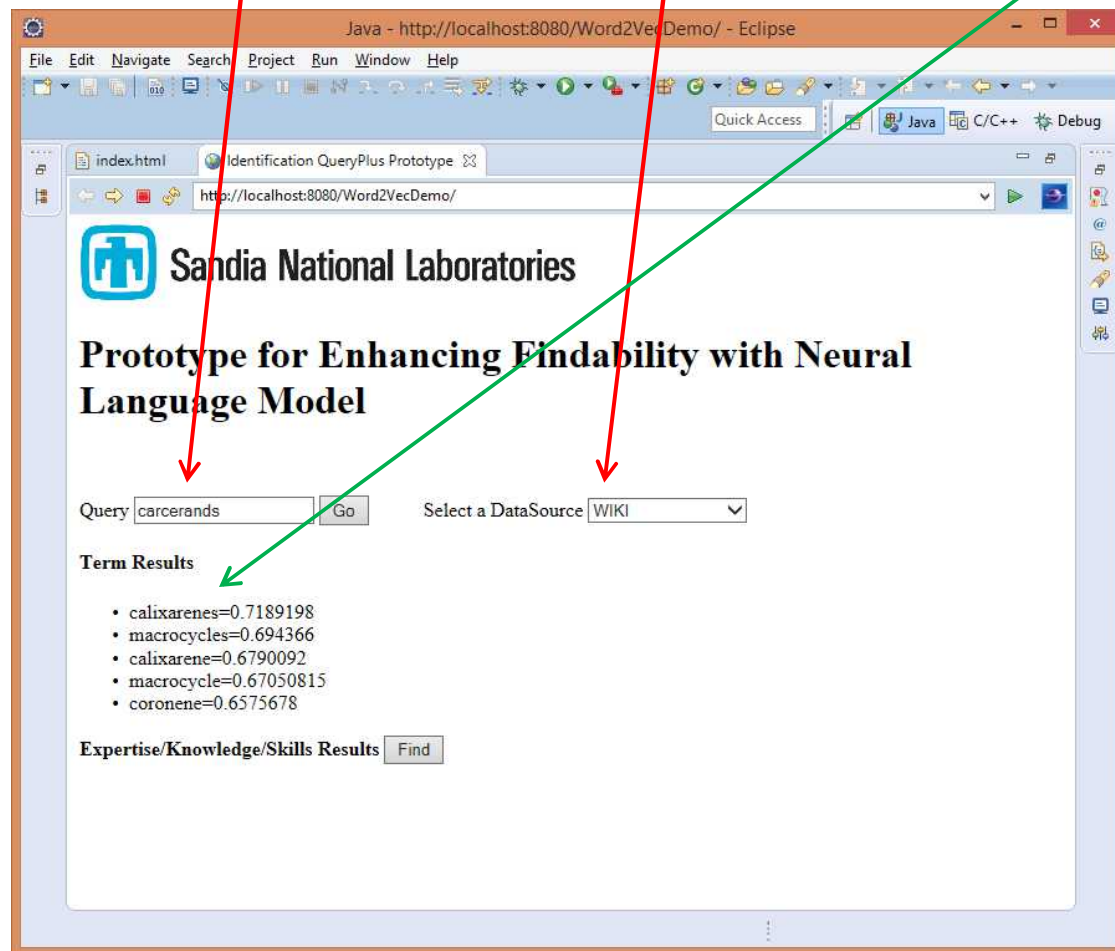


- Deliver the relevant information even the query terms not in enterprise dataset

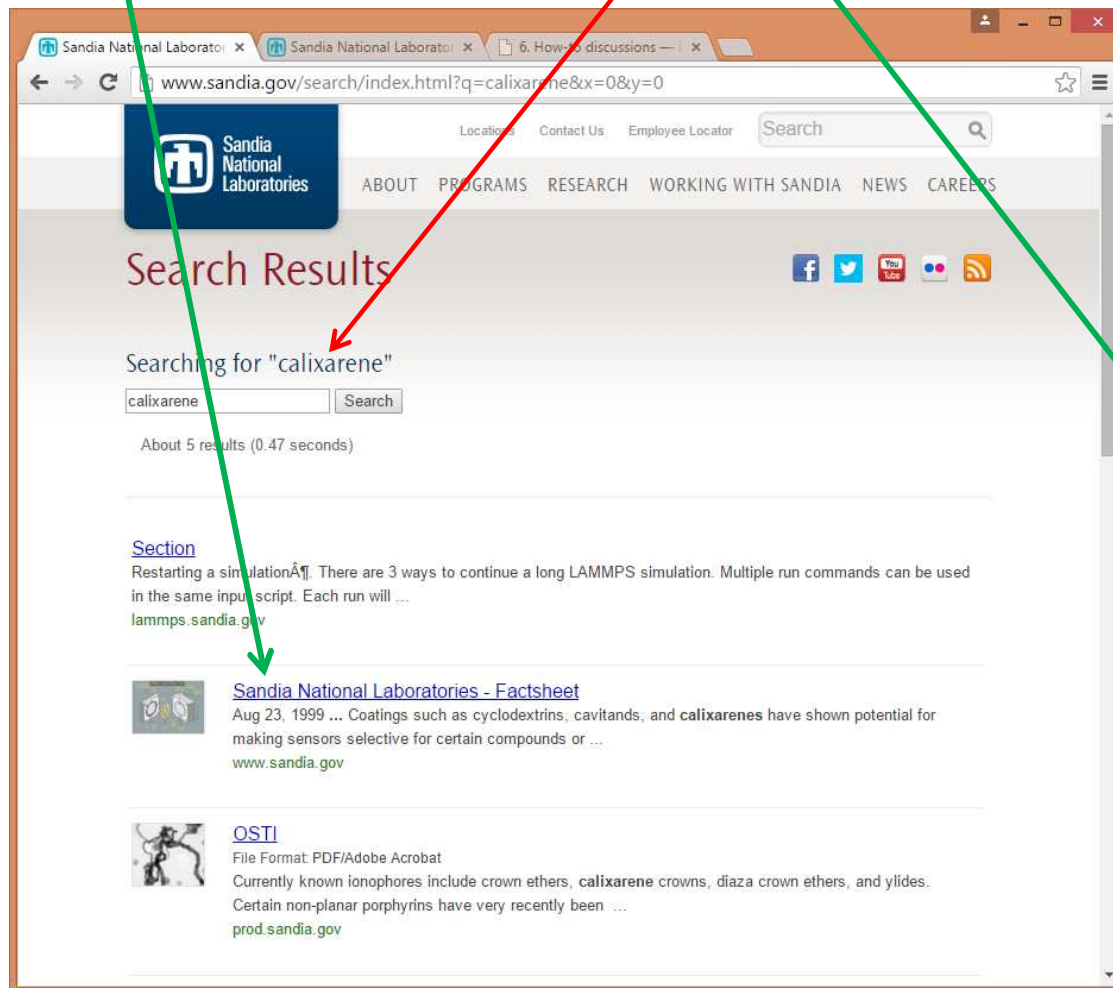
Query term "**carcerands**" not the in enterprise dataset



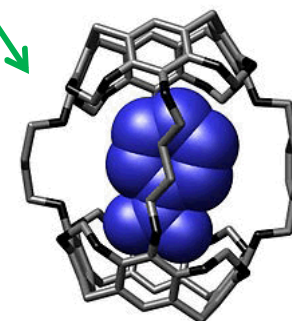
Term "carcerands" is found in Wikipedia English and a set of close terms are generated from NNLM



Sandia has a number of documents associated with "cavcerands" for which it is related to "calixarenes"



The screenshot shows a web browser window with the URL www.sandia.gov/search/index.html?q=calixarene&x=0&y=0. The page displays the Sandia National Laboratories logo and navigation links. The search results section shows "Searching for 'calixarene'" with a search bar containing "calixarene" and a "Search" button. Below the search bar, it indicates "About 5 results (0.47 seconds)". The results list includes a link to "Section" with a description about restarting a simulation, a link to "Sandia National Laboratories - Factsheet" dated Aug 23, 1999, and a link to "OSTI" with a description about ionophores. Green arrows point from the text above to the search bar and the "Sandia National Laboratories - Factsheet" link. A red arrow points from the text above to the search bar.



Wanda Sliwa and Cezary Kozlowski

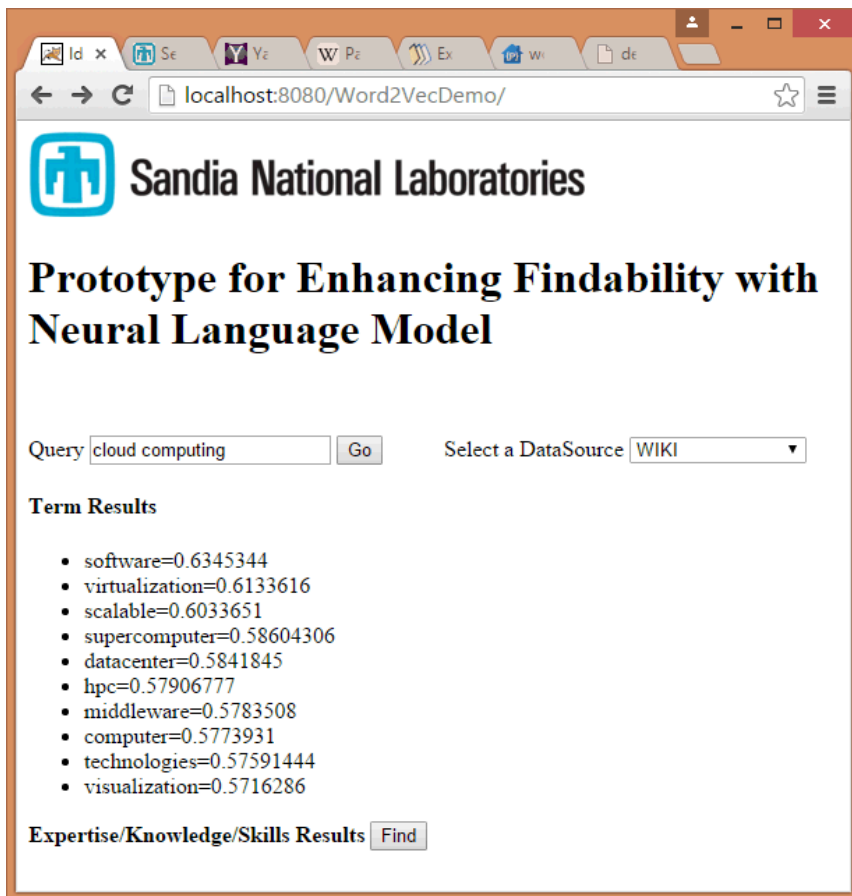
WILEY-VCH

Calixarenes and Resorcinarenes

Synthesis, Properties and Applications



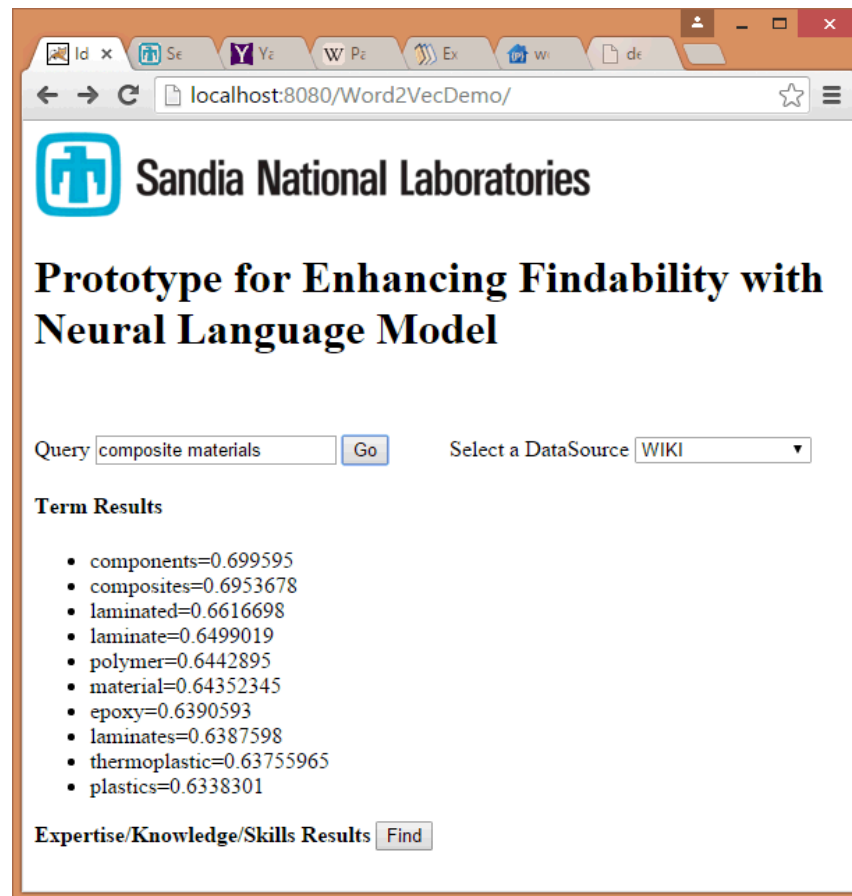
Increase the quality of retrieved information delivered: **Relevant**



The screenshot shows a web browser window with the URL `localhost:8080/Word2VecDemo/`. The page features the Sandia National Laboratories logo and the title "Prototype for Enhancing Findability with Neural Language Model". Below the title, there is a search bar with the query "cloud computing" and a "Go" button. To the right, a dropdown menu labeled "Select a DataSource" is set to "WIKI". Under the heading "Term Results", a list of terms with their corresponding scores is displayed:

- software=0.6345344
- virtualization=0.6133616
- scalable=0.6033651
- supercomputer=0.58604306
- datacenter=0.5841845
- hpc=0.57906777
- middleware=0.5783508
- computer=0.5773931
- technologies=0.57591444
- visualization=0.5716286

At the bottom, there is a section labeled "Expertise/Knowledge/Skills Results" with a "Find" button.



The screenshot shows the same web browser window as the left one, but with the query "composite materials" entered in the search bar. The "Go" button is highlighted in blue. The "Term Results" list is as follows:

- components=0.699595
- composites=0.6953678
- laminated=0.6616698
- laminate=0.6499019
- polymer=0.6442895
- material=0.64352345
- epoxy=0.6390593
- laminates=0.6387598
- thermoplastic=0.63755965
- plastics=0.6338301

The "Expertise/Knowledge/Skills Results" section and "Find" button are also present at the bottom.



Increase the quality of retrieved information delivered: **Related**

localhost:8080/Word2VecDemo/

Sandia National Laboratories

Prototype for Enhancing Findability with Neural Language Model

Query Select a DataSource

Term Results

- algorithms=0.5033293
- information=0.42070806
- datasets=0.41232347
- sparse=0.40906107
- meshes=0.40375274
- entity=0.40171173
- queries=0.40126213
- text=0.3987323
- metadata=0.3966305
- relational=0.39580885

Expertise/Knowledge/Skills Results

localhost:8080/Word2VecDemo/

Sandia National Laboratories

Prototype for Enhancing Findability with Neural Language Model

Query Select a DataSource

Term Results

- searchpoint=0.57016754
- omnifind=0.56707025
- mysites=0.55551994
- keyword=0.5526862
- application=0.5500424
- federated=0.54607517
- solr=0.54575837
- personalization=0.5381358
- searchpointnext=0.5326619
- primo=0.5207506

Expertise/Knowledge/Skills Results





Challenges of Applying NNLM in Enterprise Search

◆ Domain Specific?

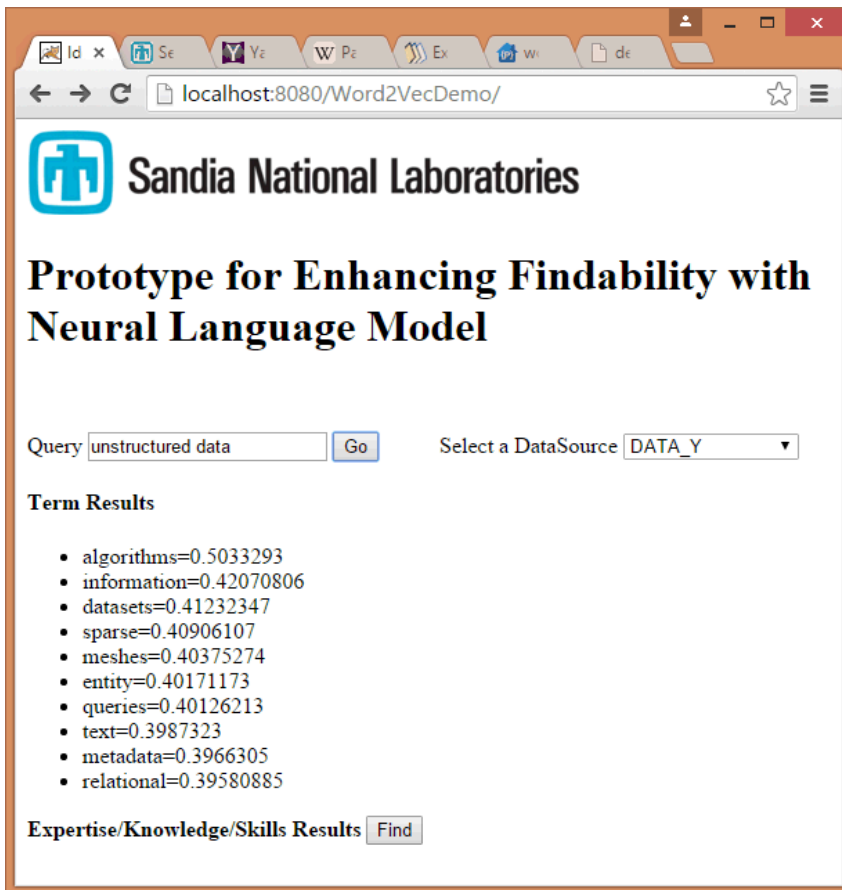
- A NNLM trained from a domain may not be directly applied to other domains in a data repository with multiple domains

◆ Maintenance?


- Data in an organization are updated dynamically, there is a need to update NNLM accordingly.



Differences among domains



localhost:8080/Word2VecDemo/

 Sandia National Laboratories

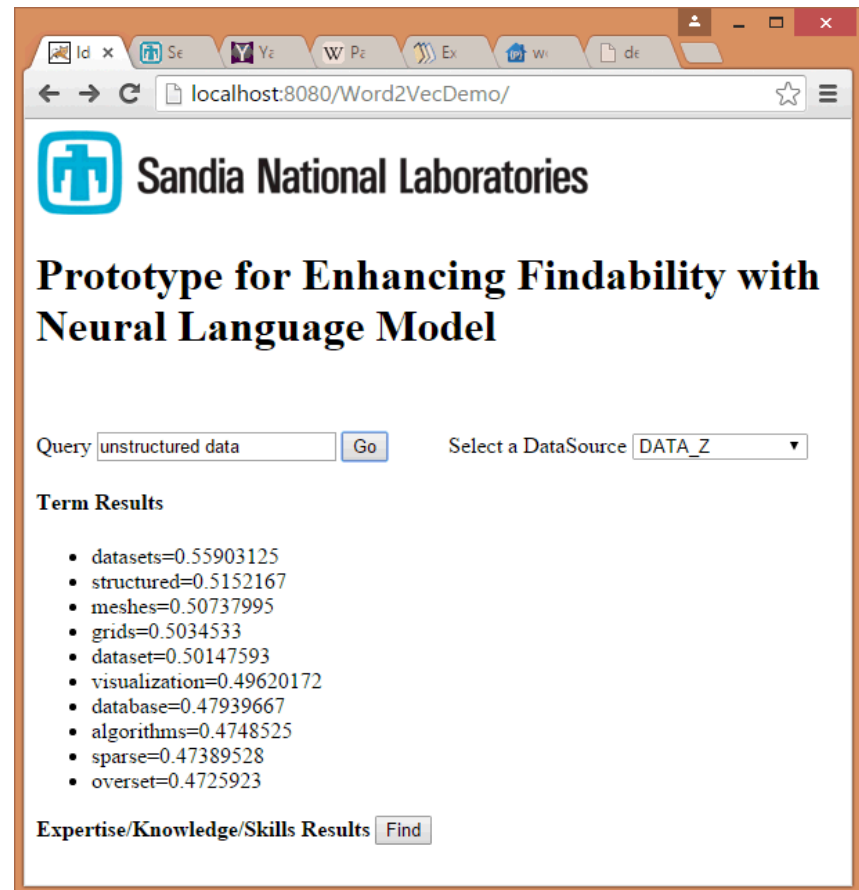
Prototype for Enhancing Findability with Neural Language Model

Query Select a DataSource


Term Results

- algorithms=0.5033293
- information=0.42070806
- datasets=0.41232347
- sparse=0.40906107
- meshes=0.40375274
- entity=0.40171173
- queries=0.40126213
- text=0.3987323
- metadata=0.3966305
- relational=0.39580885

Expertise/Knowledge/Skills Results



localhost:8080/Word2VecDemo/

 Sandia National Laboratories

Prototype for Enhancing Findability with Neural Language Model

Query Select a DataSource

Term Results

- datasets=0.55903125
- structured=0.5152167
- meshes=0.50737995
- grids=0.5034533
- dataset=0.50147593
- visualization=0.49620172
- database=0.47939667
- algorithms=0.4748525
- sparse=0.47389528
- overset=0.4725923

Expertise/Knowledge/Skills Results





Conclusions and Further Efforts

- ◆ Enterprise search can significantly improved in:
 - Deliver relevant information even the query terms not appear in the enterprise dataset
 - Increase the quality of retrieved information
 - Expand retrieved information to closely related documents
- ◆ Implementation of NNLM in Enterprise Search requires minor search engine reconfiguration
- ◆ Efforts are needed to improve training a NNLM to cover more domains in an enterprise data repository
- ◆ Further developments are needed to update NNLMs in real time

