VANDERBILT UNIVERSITY
Jean *and* Alexander Heard Libraries

# Topic Modeling:
# An Exploration

Text Analysis Community of Practice Presentation

Shenmeng Xu
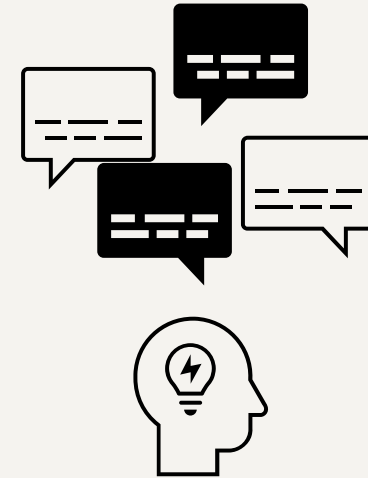Librarian for Scholarly Communications
Digital Lab
Vanderbilt Libraries

Nov 2023

# Understanding Topics…



Keyword co-occurrence map created using VOSviewer, mapping research landscape in LLM research in the field of Information and Library Science based on publication data from Web of Science
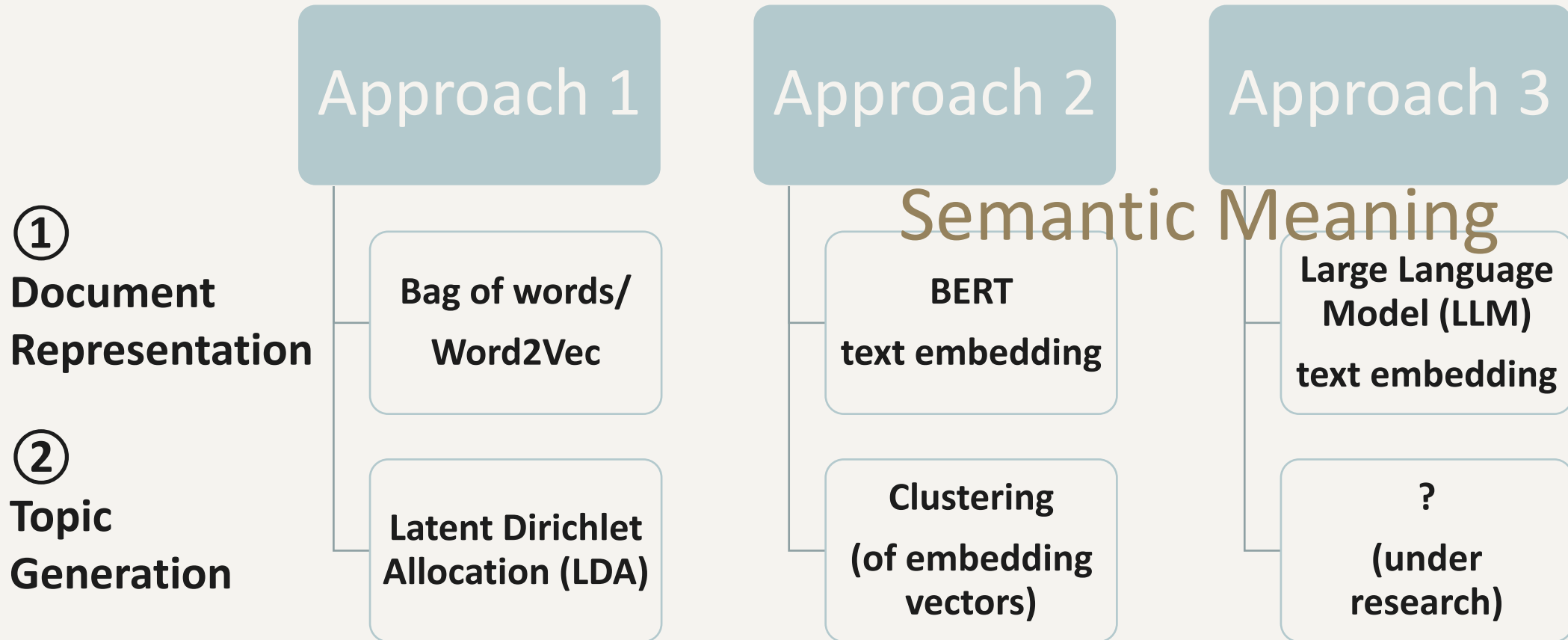
# Topic Modeling

- unsupervised machine learning
- uncovers hidden and abstract themes or topics within

    a collection of text documents
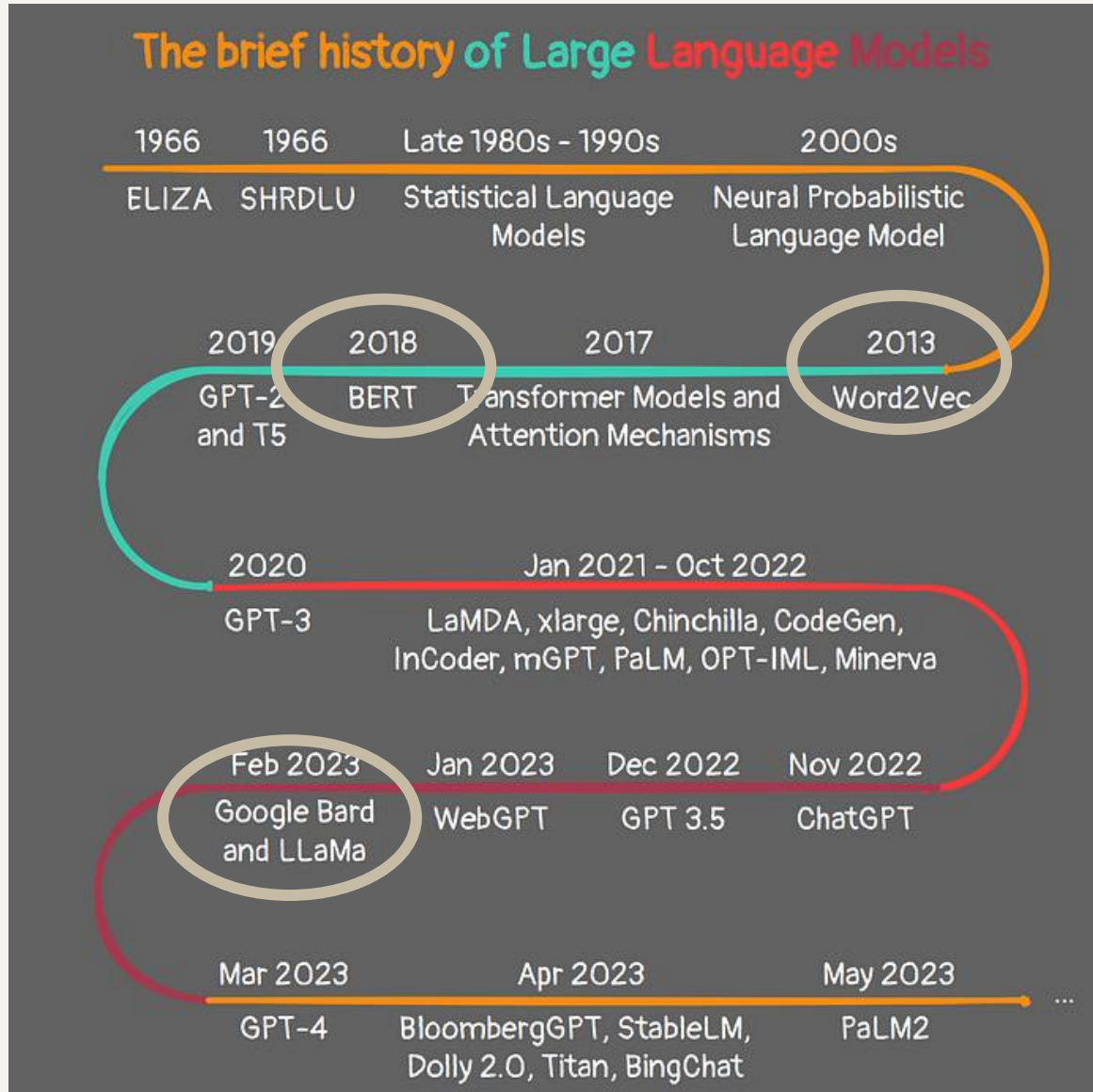- makes sense of unstructured data

# Experimenting with Different Approaches

| | Approach 1 | Approach 2 | Approach 3 |
|---|---|---|---|
| ① Document Representation | Bag of words/ Word2Vec | BERT text embedding | Large Language Model (LLM) text embedding |
| ② Topic Generation | Latent Dirichlet Allocation (LDA) | Clustering (of embedding vectors) | ? (under research) |

Semantic Meaning

# Experimenting with Different Approaches



Graph by Armin Norouzi
Reference: The brief history of Large Language Models: A Journey from ELIZA to GPT-4 and Google Bard
https://levelup.gitconnected.com/the-brief-history-of-large-language-models-a-journey-from-eliza-to-gpt-4-and-google-bard-167c614af5af

# Datasets

Dataset used in our experiments: newsgroups posts on 20 topics (split into a training set and a test set) https://huggingface.co/datasets/rungalileo/20_Newsgroups_Fixed

How to find other datasets:

- Hugging Face: https://huggingface.co/datasets

- Papers with code: https://paperswithcode.com/task/topic-models

# Colab Notebook

https://colab.research.google.com/drive/18YM5HoLCi1_gb39OaiCgVnaRejAUyN4m?usp=sharing

# Additional Resources

- Word Embeddings: https://ai.engin.umich.edu/2018/07/23/word-embeddings-and-how-they-vary/

- LDA: https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation

- LDA Evaluation: https://towardsdatascience.com/evaluate-topic-model-in-python-latent-dirichlet-allocation-lda-7d57484bb5d0

- BERTopic: https://maartengr.github.io/BERTopic/index.html

- The Illustrated Transformer: https://jalammar.github.io/illustrated-transformer/

Thank you!