# Pretraining Language Models for Synonymy Prediction at Scale in the UMLS Metathesaurus

Thilini Wijesiriwardene

Mentors:
Dr. Vinh Nguyen
Dr. Olivier Bodenreider

# Motivation

- UMLS Metathesaurus integrates biomedical terms from various vocabularies
- Current UMLS construction process: tedious, error-prone, expensive
- Our prior work*:
  - LexLM: a deep learning model leveraging lexical patterns
  - Rule-based approximation of current construction process
- Can more recent techniques in Deep Learning and NLP perform better in UMLS Metathesaurus construction?

# Objectives

Develop UBERT, a BERT based language model pretrained on UMLS data and Synonymy Prediction task that can provide state-of-the-art performance on Synonymy Prediction

*Nguyen, V., Yip, H. Y., & Bodenreider, O. (2021, April). Biomedical Vocabulary Alignment at Scale in the UMLS Metathesaurus. In *Proceedings of the Web Conference 2021* (pp. 2672-2683)

# Synonymy Prediction Task*

"atom" – Single term form a source vocabulary

atoms with same meaning are grouped in to one concept identified by a Concept Unique Identifier (CUI)
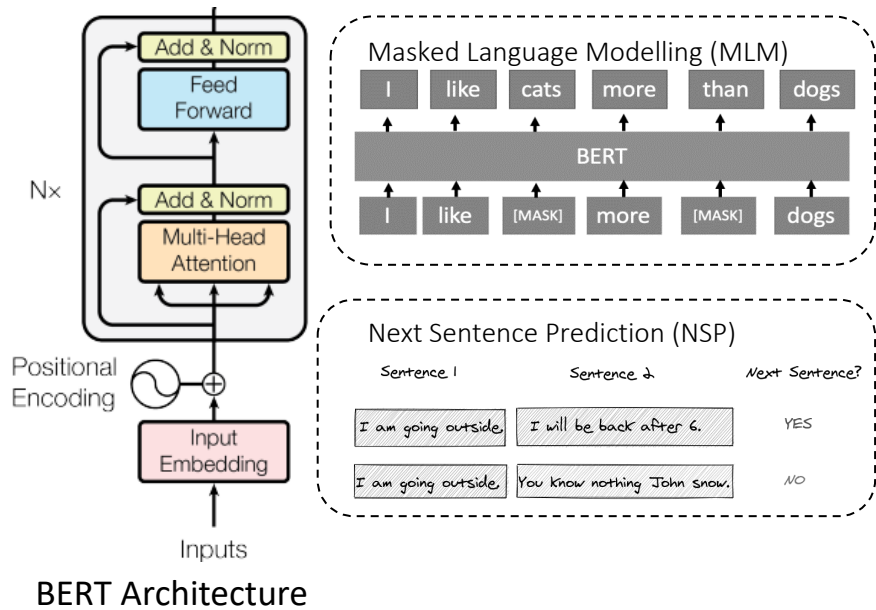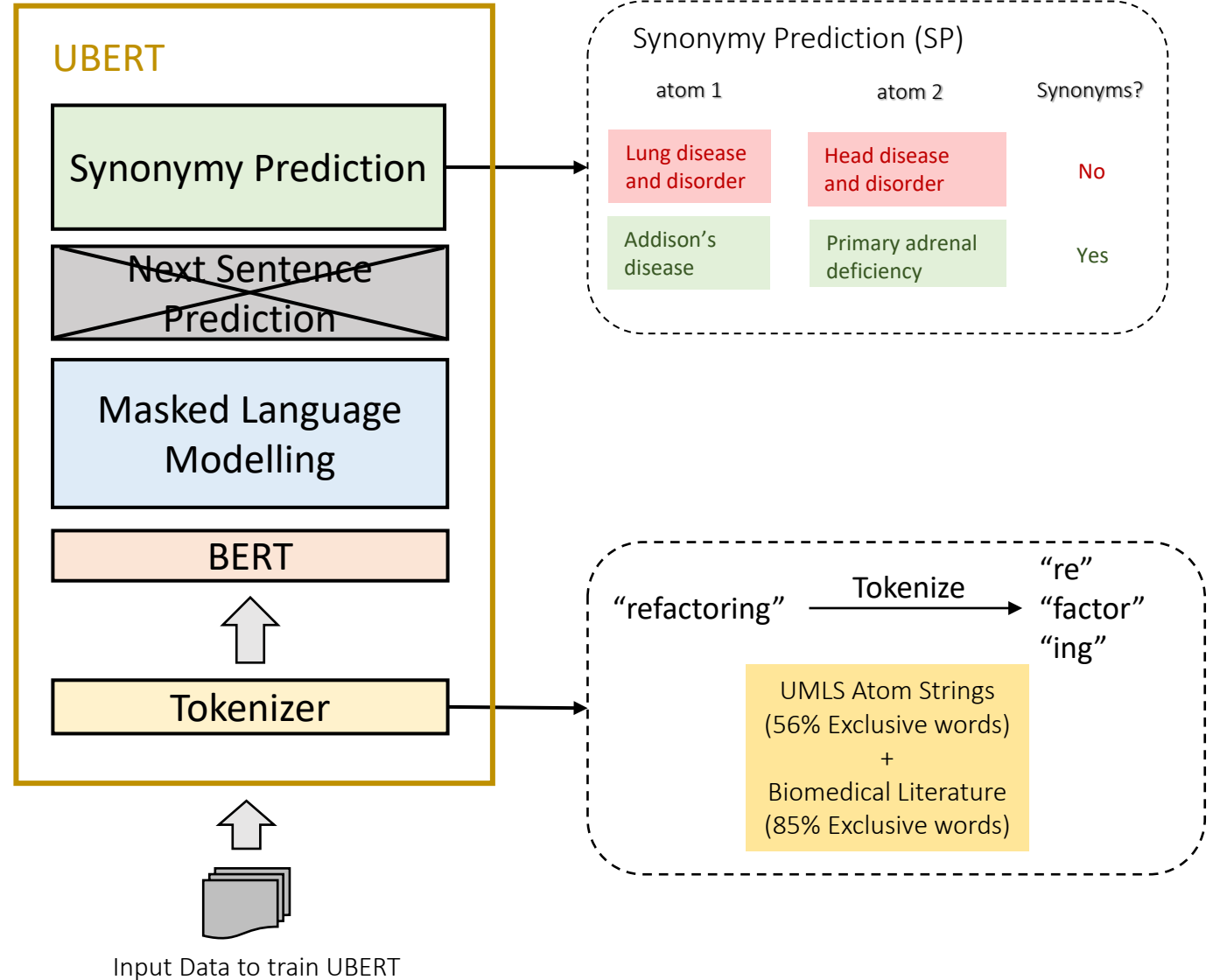


*Nguyen, V., Yip, H. Y., & Bodenreider, O. (2021, April). Biomedical Vocabulary Alignment at Scale in the UMLS Metathesaurus. In *Proceedings of the Web Conference 2021* (pp. 2672-2683).

NIH National Library of Medicine
Lister Hill National Center for Biomedical Communications

# BERT (Bidirectional Encoder Representations from Transformers)*

Artificial neural network-based language model, designed to provide meaning for a word by using its surrounding context.
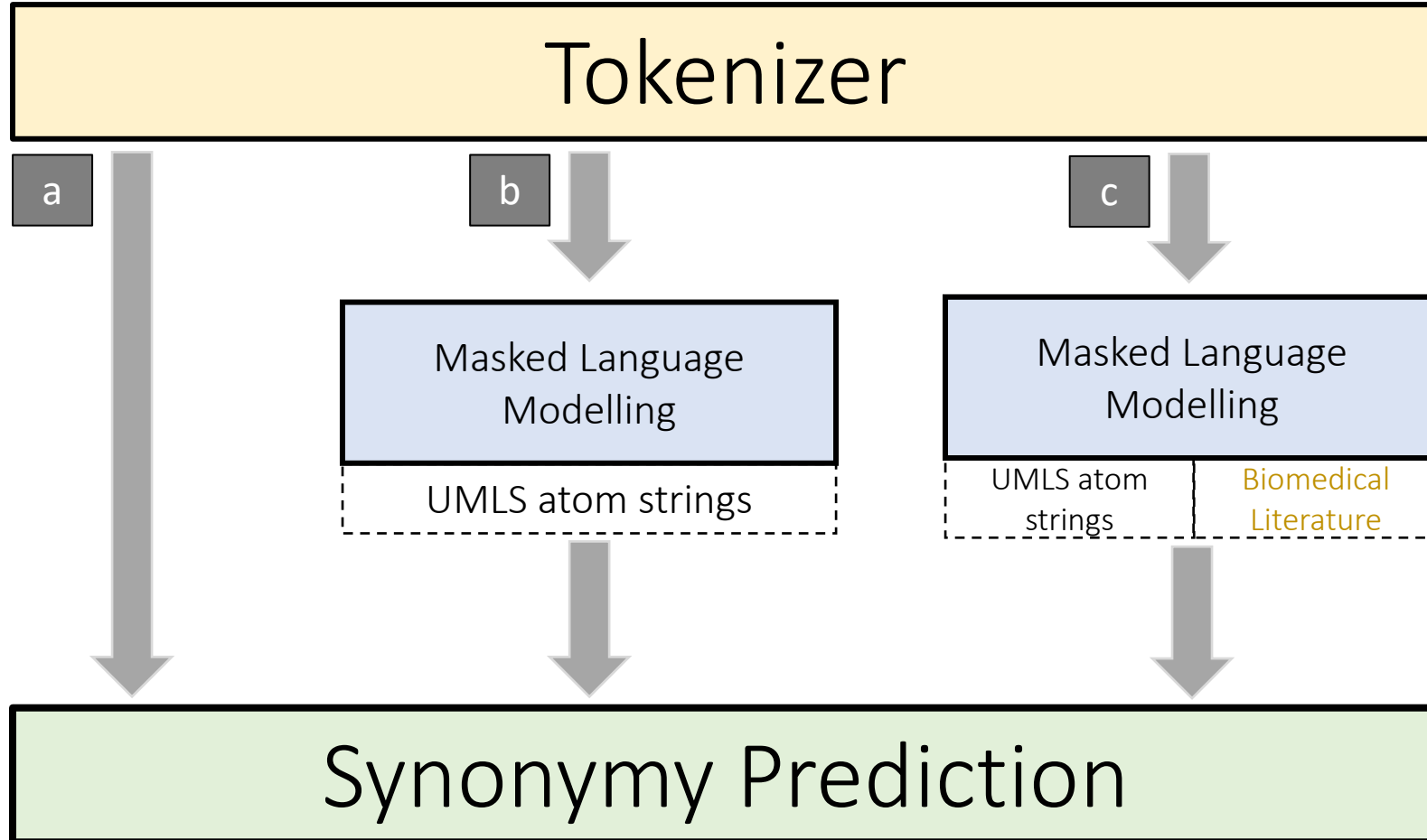


BERT Architecture

*Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019, January). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In NAACL-HLT (1).

# UBERT Architecture



Input Data to train UBERT
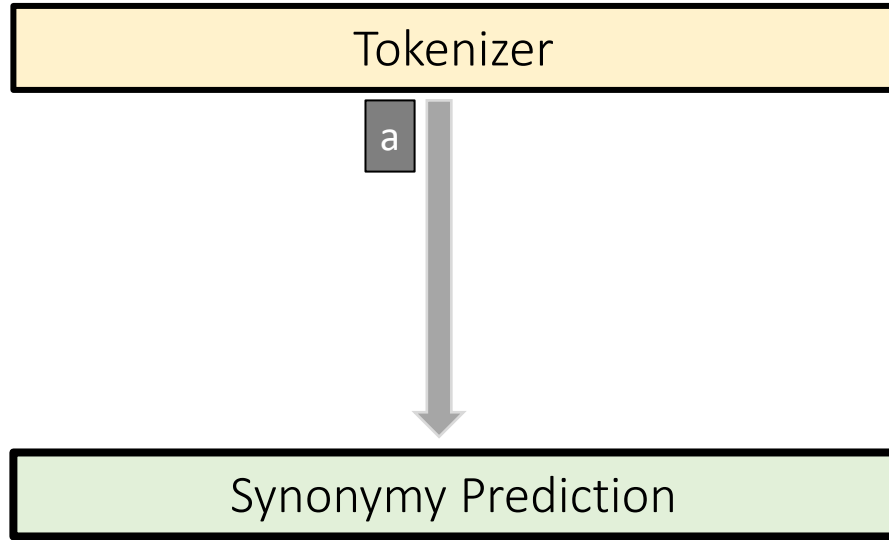
# Datasets for Training & Testing UBERT

- Train tokenizer and Masked Language Modelling task
  - PubMed Abstracts ➔ ~ 4.5 billion words
  - PubMed Central Full Texts ➔ ~ 13.5 billion words
- Synonymy Prediction task :
  - Training ➔ 118 million pairs
    - Positive (synonymous) pairs
    - *Negative (non-synonymous) pairs - varying degrees of lexical similarity
  - Testing ➔ 170 million pairs

*Nguyen, V., Yip, H. Y., & Bodenreider, O. (2021, April). Biomedical Vocabulary Alignment at Scale in the UMLS Metathesaurus. In *Proceedings of the Web Conference 2021* (pp. 2672-2683).

# Training Setup

# Initial Results



| Metric | RBA | LexLM | UBERT |
|---|---|---|---|
| F1 | 0.7651 | 0.9061 | **0.9974** |
| Accuracy | 0.9863 | 0.9938 | **0.9950** |
| Precision | 0.8631 | 0.8875 | **0.9991** |
| Recall | 0.6871 | 0.9254 | **0.9957** |

Training UBERT for a single epoch takes around ~3 hours on 16 Nvidia V100X GPUs.

# Internship Progress

Develop UBERT, a BERT based language model pretrained on UMLS data and Synonymy Prediction task that can provide state-of-the-art performance on

1. Synonymy prediction
   - Variant a
   - Variant b
   - Variant c

# Conclusions and Future Work

- UBERT is a potential candidate for UMLS Metathesaurus construction

- Look for more efficient architectures for Synonymy Prediction task that leverage atoms' contexts when predicting synonymy.

- Evaluating performance of UBERT in BioNLP tasks such as biomedical Named Entity Recognition, biomedical Relations Extraction, etc.

# Acknowledgements

Dr. Olivier Bodenreider

Dr. Vinh Nguyen

Goonmeet

Vishesh

Joey

Dr. Kin Wah Fung

Dr. Yuqing Mao

# Thank You
# Questions?

# References

1. Bodenreider, O. (2004). The unified medical language system (UMLS): integrating biomedical terminology. 3 Nucleic acids research, 32(suppl_1), D267-D270.
2. https://medium.com/genei-technology/richer-sentence-embeddings-using-sentence-bert-part-i-ce1d9e0b1343
3. https://amitness.com/2020/02/albert-visual-summary/
4. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019, January). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In NAACL-HLT (1).
5. Nguyen, V., Yip, H. Y., & Bodenreider, O. (2021, April). Biomedical Vocabulary Alignment at Scale in the UMLS Metathesaurus. In *Proceedings of the Web Conference 2021* (pp. 2672-2683)
6. https://huggingface.co/transformers/model_doc/bert.html