

# Internet-NLP: Allowing Natural Language Processing Models to Connect to the Internet

**Thamognya Kodi**

Ruamrudee International School

[contact@thamognya.com](mailto:contact@thamognya.com)

## Abstract

In this paper, I present **Internet-NLP**, a new control-flow wrapper abstraction to enable the utilization of data from the internet (or a knowledge-database when offline) for existing context-needing Natural Language Processing (NLP) models to function without any given context. Internet-NLP can be used, finetuned alongside existing NLP models via its config settings and additionally its Long Short Term Memory neural network (LSTM neural network) can also be trained. Additionally incorporations of Masked Language Models (MLM) such as BERT, or LinkBERT (Devlin et al., 2019; Yasunaga et al., 2022a) can be utilized to improve search queries, and therefore retrieve more accurate and reliable data. Furthermore, Internet-NLP utilizes a LSTM, Reinforcement Learning and caches to allow for multi-turn NLP tasks, and improvement via Reinforcement Learning from user.

Internet-NLP, in basic terms, provides the context for context-needing NLP models to let them function. Internet-NLP can be improved via finetuning, and training of LSTM and Reinforcement Learning model (which can be trained alongside the NLP model), which enables for better search queries, and subsequently results. It obtains state-of-the-art (SOTA) results in QA and NLI without context.

Additionally in this paper, I also present new NLP and Natural Language Inference (NLI) models to assist **Internet-NLP**:

- Open-book question and long answer (QA) via GPT-NeoX-20b (Andonian et al., 2021; Black et al., 2022)
- CrossEncoder NLI via LinkBERT (Reimers and Gurevych, 2019a; Thakur et al., 2021; Yasunaga et al., 2022a)
- Answer to context NLP via T5 (Raffel et al., 2019a)

Along with these models, I also present new general purpose QA and NLI datasets:

- ALoT NLI made from datasets: ANLI, MultiNLI, and SNLI (Nie et al., 2020; Williams et al., 2018; Bowman et al., 2015)
- ALoT OpenBookQA made from datasets: CoQA, Natural Questions, and SQuAD (Reddy et al., 2018; Kwiatkowski et al., 2019; Rajpurkar et al., 2018)

As a result of these Internet-NLP, models and datasets the accuracy and reliability of most context-needing NLP models on most NLP tasks, especially tasks that require more factual responses with no given context increased.

Internet-NLP and the new NLP and NLI models, which were trained on the general-purpose datasets (ALoT NLI, and ALoT OpenBookQA). Internet-NLP, by default utilizes an Text-Generative model GPT-NeoX (Andonian et al., 2021; Black et al., 2022) for long responses and LinkBERT (Yasunaga et al., 2022a) for short responses. For 2 choices (for ex: True and False) Bi-Encoder NLI has been used and for multiple choices CrossEncoder will be used (Thakur et al., 2021).

## 1 Introduction

There are currently two main solutions for utilizing NLP tasks with no context provided:

1. Large Pre-Trained Text-Generation and Text2Text-Generation Model
  - Pre-trained Text-generation models, like GPT-NeoX, GPT-3, and etc. (Black et al., 2022; Andonian et al., 2021; Brown et al., 2020) can be trained for open-domain question-answering closed-book language model tasks (ODQA LM) (Weng, 2020). When used in ODQA tasks, they achieve SOTA results in such tasks, have high accuracy and are fast but are much larger in size than open-book (context-needing) language models.
  - Additionally Pre-trained Text2Text-generation models, like T5 (Raffel et al.,

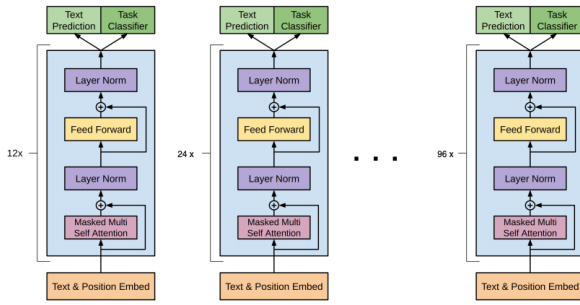


Figure 1: This is an illustration of the architecture of GPT-2 and GPT-3, a popular Text-Generation model (Dzlab; the).

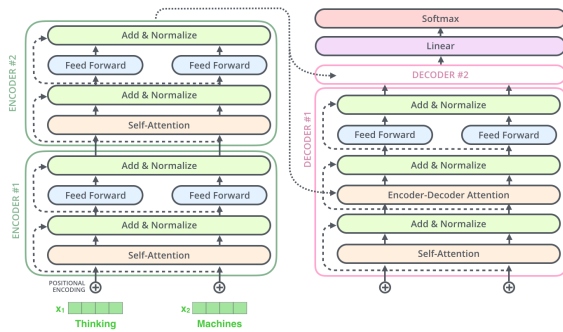


Figure 2: This is an illustration of the architecture of T5, a popular Text2Text-Generation model (Alammar).

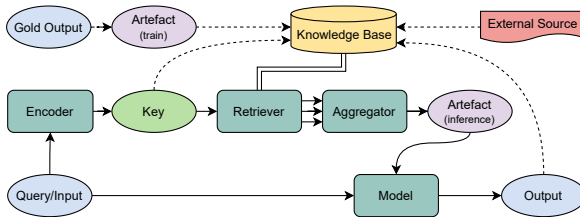


Figure 3: This is an illustration of how LMs with a knowledge base and artifact retriever work (Zouhar et al., 2022).

2019a) that have open-domain question-answering closed-book (no context) language models (ODQA LM) capabilities (Weng, 2020). These closed-book QDQA LMs are comparatively state-of-the-art performance in many no-context NLP tasks, mainly question-answering. Text2Text-generation models for such no-context NLP tasks are usually large, slow, and have a low accuracy (Roberts et al., 2020a).

- Example: T5 (Raffel et al., 2019a)
- Illustration of how ODQA LM work: 1

## 2. Large Knowledge Database with a Context-Needing Language Model

- Large Knowledge base with an pre-trained

open-book LM and retriever, provides a comparatively higher performance, accuracy and the model itself is small. These models however require — usually — a large knowledge base which makes the overall solution large, but still fast and with more accuracy on the field the knowledge base specializes in.

- Example: LinkBERT (Yasunaga et al., 2022b) with an artifact retriever (Zouhar et al., 2022) with a knowledge base such as DBpedia or WikiData (Auer et al., 2007; Vrandečić and Krötzsch, 2014)
- Illustration of how LMs with a knowledge base and artifact retriever work: 3

Solution 1 and 2 achieve the same end goal of NLP tasks without context via two different methods; these current solutions restrict NLP tasks and accuracy without context, especially for more open-domain tasks. The major limitation in this case would be accuracy, efficiency and size of models and their knowledge base which then limit the use cases of closed-book open-domain NLP tasks.

In this paper, I propose Internet-NLP, an direct improvement to solution 1 which allows NLP models to not require a large knowledge base (although you can configure Internet-NLP to utilize a knowledge base) that incorporates the internet's vast knowledge along utilizing data in hyperlinks in webpages (Yasunaga et al., 2022b) to create a more resource-filled data for our existing or future context-needing pre-trained model to use for NLP tasks. Internet-NLP encompasses pre-trained NLP and NLI models, along with its web-data-scraper creates an small temporary on-basis data and a cache for NLP tasks to be performed without given context.

Utilizing the vast data on the internet, graph of documents as corpus (Yasunaga et al., 2022b) allows us to enable to reduce our solution size, increase efficiency and increase accuracy. Additionally unlike usage of static data, Internet-NLP utilizes the dynamic, and frequent updating data of the internet which enables us to utilize any type of NLP model along with NLI models to allow us to follow a sequence of control flow to get the context for context-needing models. This approach utilizes a combination of data-collection (Hrkút et al., 2020) for NLPs with context-needing open-domain NLP to gain more accurate results in most no-context NLP tasks.

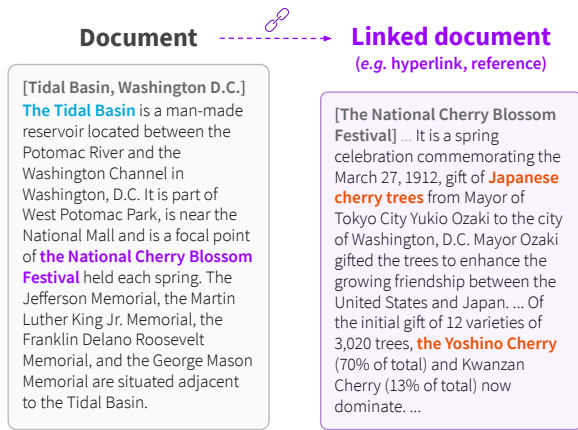


Figure 4: This is an illustration of example of how LinkBERT utilizes hyperlinks to make a graph corpus (Yasunaga et al., 2022b).

Additionally Internet-NLP’s Text2Text-generation search query model: T5 (Raffel et al., 2019a) and LSTM noun remembrance using parts of speech tagging (Chiche and Yitagesu, 2022) on ALotClosedBookQA with it improving search queries based on the difference on answer recieved and the answer from datasets, using parts of speech tagging on answers (Chiche and Yitagesu, 2022).

## 2 Related Work

### 2.1 Internet-NLP

#### 2.1.1 NLP Models with Knowledge Base and Retriever

These approaches are one the two most popular current solution for NLP tasks to be done without context. It utilizes an knowledge base, a retriever for this data and a LM depending on the use case for example (this list is not extensive):

- question answering: LinkBERT or T5 (Yasunaga et al., 2022b; Raffel et al., 2019a)
- NLI: CrossEncoder models BERT or DeBERTa (Thakur et al., 2021; Devlin et al., 2018; He et al., 2020)

This allows for no-context NLP applications (especially question and answering) to function without any context given, due to knowledge base and retriever providing the context. An representation of this is shown in illustration 3 (Zouhar et al., 2022).

### 2.2 Internet-NLP’s NLP models

#### 2.2.1 LinkBERT

LinkBERT is a NLP model that is a pre-trained BERT (Devlin et al., 2018) model that is trained on

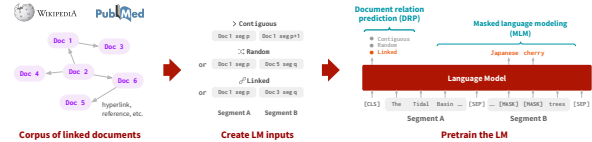


Figure 5: This is an illustration of example of how LinkBERT makes a graph corpus (Yasunaga et al., 2022b).

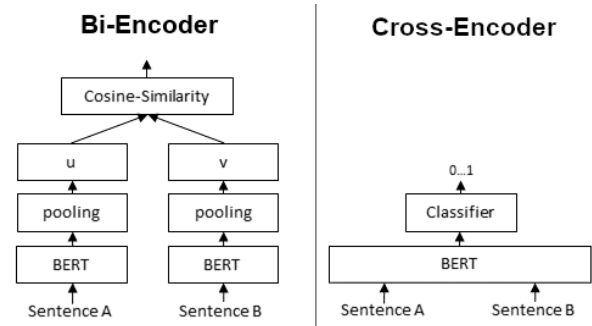


Figure 6: This is an illustration of how NLI using Cross-Encoders vs Bi-Encoder work like (Thakur et al., 2021).

a graph-based corpus of documents from not only documents but also the hyperlinks in documents. It utilizes a "fusion of graph-based and language-based self-supervised learning" (Yasunaga et al., 2022b). It gains better performance on graph-based data corpus than other pre-trained NLP models due to it being trained with utilizing graph-based self-supervised learning.

These are illustrations that explain LinkBERT’s graph-based and language-based fusion:

- This illustration shows how hyperlinks can contain crucial information: 4.
- This illustration shows how LinkBERT (Yasunaga et al., 2022b) makes a graph from links: 5.

For training the Internet-NLP and LM for Text2Text-generation for question answering would be utilizing the fusion of graph-based and language-based learning LinkBERT revolutionized (Yasunaga et al., 2022b).

### 2.3 Internet-NLP’s NLI models

#### 2.3.1 Cross-Encoder NLI Models

NLI compares two sentences to given an output of entailment (true), neutral or contradiction (false).

Utilizing Cross-Encoder for NLI applications that allow for the utilization of Cross-Encoder (an illustration of Cross-Encoders 6) where two sentence are passed simultaneously, and then utilizing

a classifier to get the output of 0 to 1 which goes from contradiction to entailment (Thakur et al., 2021; Reimers and Gurevych, 2019b).

### 3 Preliminaries

The preliminaries listed are NLP tasks Internet-NLP benefits from the access to internet listed:

#### 3.1 Question Answering

The tasks of training an NLP model to utilize question, and context (or in the case of ODQA closed-book LM just question) to create a logical answer. The current most popular would be context-needing question answering models, where in the answer is provided in the context. These models utilize reading comprehension to utilize the context to make an answer based on the question (Roberts et al., 2020b).

Closed-book QDQA LM are a type of question-answering model where there is no context provided, and these are usually the hardest variant to train, and results in large sizes, low efficiency, and low accuracy. These models can only be asked context-independent questions such as facts (Roberts et al., 2020b).

The alternative to ODQA LM would be utilizing a knowledge base and retriever for getting the required context from a knowledge base and then utilizing an context-needing question-answering NLP model which would be known as open-book question-answering model (Roberts et al., 2020b). This however require knowledge base which would be static and hence would not contain the latest information; additionally requires a large database and hence the however solution is also large.

In this publication, Internet-NLP applies question answering open-book LM with the constraint of not utilizing a knowledge base and keeping the overall solution size to be low, high efficiency and high accuracy with the ability to also being asked context-dependent (by giving the optional context) and context-independent questions. Internet-NLP utilizes the internet to replace the knowledge base, utilize a retriever to get the required information from the internet data and then an open-book Text2Text-generation model to create an answer from the information, question and any extra context given.

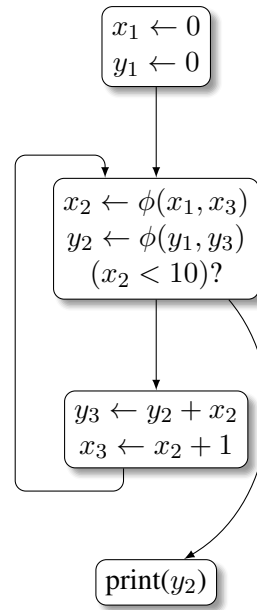


Figure 7: This is an illustration of how Internet-NLP's control flow works.

#### 3.2 Natural Language Inference

NLI models require a premise (similar to context) and hypothesis (an prediction) to give one of the following: entailment (hypothesis is correct based on premise), neutral (hypothesis is neither correct nor wrong based on premise) and contradiction (hypothesis is wrong based on premise).

Current no-premise NLI models utilizes a knowledge base to reproduce the premise via a retriever and then utilize an NLI model to then given output.

In this publication, Internet-NLP produces the premise based on the hypothesis by converting the hypothesis into an search query (via an Text2Text-generation LM) which will then be scraped for results and then be indivisually compared to the hypothesis to to only select ones that have either contradiction or entailment to then give an ouput on wether its an entailment or contradiction. This allows for hypothesis to be checked if they are either correct or wrong without an large knowledge base or model.

### 4 Internet-NLP

This publication will introduce Internet-NLP and its control flow for allowing NLPs to connect to internet, which will replace traditional knowledge bases with the resources on the internet.

In the control flow diagram 7, it shows how Internet-NLP gains its data for NLP tasks and also makes sure that the data scraped is accurate and not

offensive for the NLP task it is being asked to do; Internet-NLP does this by utilizing several different NLP and NLI models in combination to enable this data collection system. This allows other NLP models to utilize the data to allow for other NLP tasks that was requested.

Internet-NLP's control flow diagram 7 will be explained in the following subsections.

#### 4.1 NLP Tasks Applicable

Internet-NLP currently allow for the following NLP tasks without context:

- Question Answering
- Zero-Shot Classification
- Natural Language Inference
- Text2Text Generation
- Conversational (this still in beta and does not completely work)

#### 4.2 Disclaimers

##### 4.2.1 Types of English

Internet-NLP at this point of time can only fully understand "formal" English (For). Additionally idioms, similes, and other figures of speech are not understood by Internet-NLP or it's models.

##### 4.2.2 Output of Internet-NLP

The accuracy of the output of Internet-NLP depends on the data it scrapes which may not be completely accurate (which the chances are minimized to an extent with utilizing mutiple resources) and may contain profanity or abrasive language which may or may not affect the output.

#### 4.3 Common Components of Internet-NLP's Process

##### 4.3.1 Answer To Question Text2Text-generator

##### 4.3.2 Search Queries Text2Text-generator

The search query generator that enables converting questions into viable search queries utilizes a fastT5 model (Raffel et al., 2019b). It is trained on reddit and quora questions (that are non-mathematical i.e does not require logical computation) and then passed through an parts of speech tagging model and normalizer wherein the question is optimized for search engines by removing specific details and punctuation (Bet).

The reason for utilizing fastT5 models rather than the parts of speech tagging model comes down due to efficiency issues as fastT5 outperforms the parts of speech tagging model (Banga and Mehndiratta, 2017; Raffel et al., 2019b).

##### 4.3.3 Data Collection

#### 4.4 Question Answering

##### 4.4.1 Answer to Question Text2Text-generator

In the case of question answering without context, Internet-NLP only needs one of following:

- Question
  - In this case Internet-NLP passes the question through the Search Query Text2Text-generator 4.3.2 wherein an output of a optimized search question for search engine is returned. This optimized question will be used for data collection 4.3.3.
- Answer
  - In this case the Answer to Question Text2Text-generator 4.3.1. After which it follows the same process of question optimization explained above in Question case 4.4.1.

##### 4.4.2 Natural Language Inference Without Premise

#### References

[Formal and informal style.](#)

[How deep is the machine?](#)

[Refine web searches.](#)

Jay Alammar. [The illustrated transformer.](#)

Alex Andonian, Quentin Anthony, Stella Biderman, Sid Black, Preetham Gali, Leo Gao, Eric Hallahan, Josh Levy-Kramer, Connor Leahy, Lucas Nestler, Kip Parker, Michael Pieler, Shivanshu Purohit, Tri Songz, Wang Phil, and Samuel Weinbach. 2021. [GPT-NeoX: Large Scale Autoregressive Language Modeling in PyTorch.](#)

Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. Dbpedia: A nucleus for a web of open data. In *Proceedings of the 6th International The Semantic Web and 2nd Asian Conference on Asian Semantic Web Conference, ISWC'07/ASWC'07*, page 722–735, Berlin, Heidelberg. Springer-Verlag.

Ritu Banga and Pulkit Mehndiratta. 2017. [Tagging efficiency analysis on part of speech taggers.](#) pages 264–267.

- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. [GPT-NeoX-20B: An open-source autoregressive language model](#). In *Proceedings of the ACL Workshop on Challenges & Perspectives in Creating Large Language Models*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). *CoRR*, abs/1508.05326.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). *CoRR*, abs/2005.14165.
- Alebachew Chiche and Betselot Yitagesu. 2022. [Part of speech tagging: a systematic review of deep learning and machine learning approaches](#). *Journal of Big Data*, 9(1):10.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dzlab. [[link](#)].
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. [Deberta: Decoding-enhanced bert with disentangled attention](#).
- Patrik Hrkút, Štefan Toth, Michal Ďuračík, Matej Meško, Emil Krsak, and Miroslava Mikušová. 2020. [Data Collection for Natural Language Processing Systems](#), pages 60–70.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. [Adversarial NLI: A new benchmark for natural language understanding](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019a. [Exploring the limits of transfer learning with a unified text-to-text transformer](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019b. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *arXiv e-prints*.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don't know: Unanswerable questions for squad](#). *CoRR*, abs/1806.03822.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2018. [Coqa: A conversational question answering challenge](#). *CoRR*, abs/1808.07042.
- Nils Reimers and Iryna Gurevych. 2019a. [Sentencebert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019b. [Sentencebert: Sentence embeddings using siamese bert-networks](#).
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020a. [How much knowledge can you pack into the parameters of a language model?](#) *CoRR*, abs/2002.08910.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020b. [How much knowledge can you pack into the parameters of a language model?](#)
- Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. 2021. [Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 296–310, Online. Association for Computational Linguistics.
- Denny Vrandečić and Markus Krötzsch. 2014. [Wiki-data: A free collaborative knowledgebase](#). *Commun. ACM*, 57(10):78–85.
- Lilian Weng. 2020. [How to build an open-domain question answering system?](#) *lilianweng.github.io*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022a. [LinkBERT: Pretraining language models with document links](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8003–8016, Dublin, Ireland. Association for Computational Linguistics.

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022b. [Linkbert: Pretraining language models with document links](#).

Vilém Zouhar, Marius Mosbach, Debanjali Biswas, and Dietrich Klakow. 2022. [Artefact retrieval: Overview of nlp models with knowledge base access](#).