

# Retrieve and Refine System for Polish Question Answering

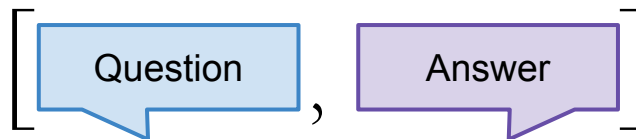
Piotr Rybak, ML Research at Allegro.pl

# Datasets

- PolEval
  - Official data provided by task organizers
  - 1000 dev + 2500 test-A
- Jeden z dziesięciu
  - 1004 deduplicated questions
- Multi-lingual Knowledge

## Questions & Answers

- Available here: [hf.co/datasets/mkqa](https://huggingface.co/datasets/mkqa)
- 10k raw questions → 1875 after filtration

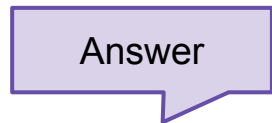


# Architecture

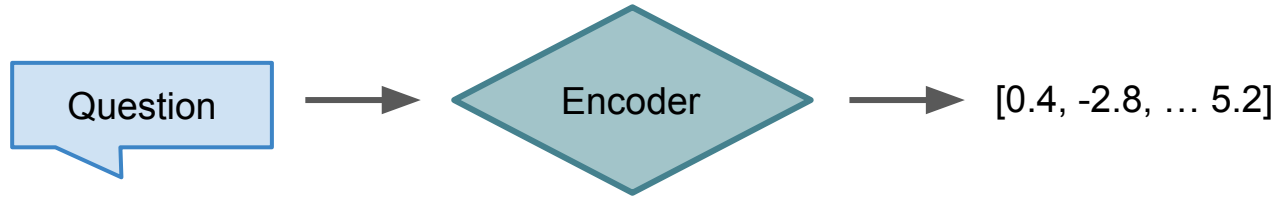
Question



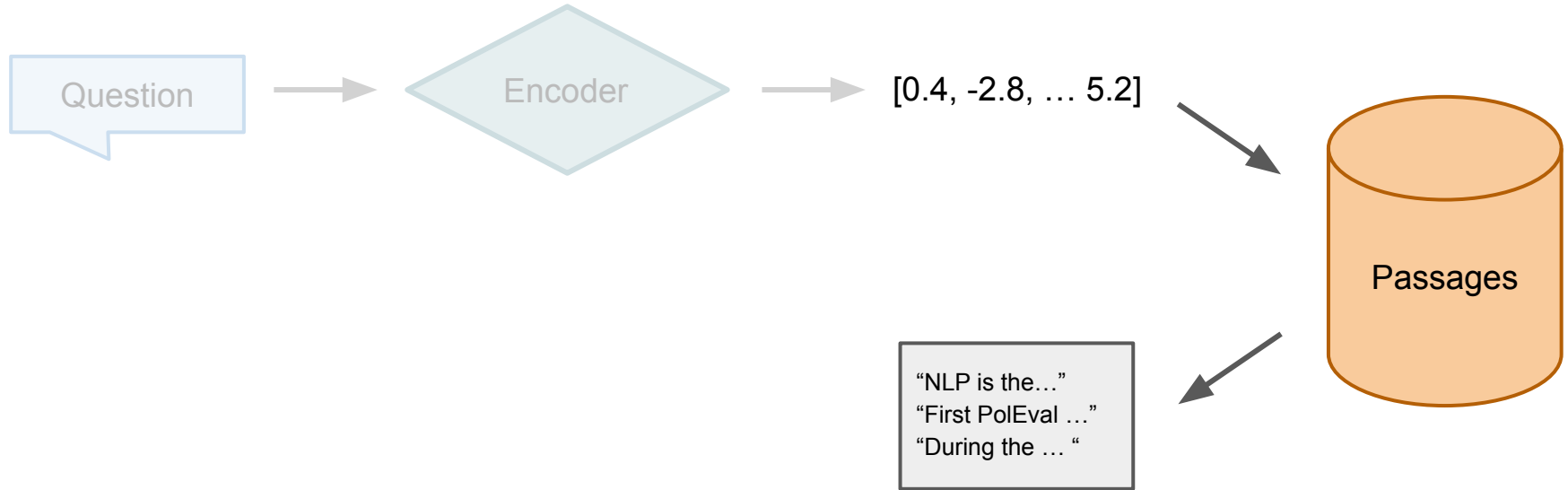
Answer



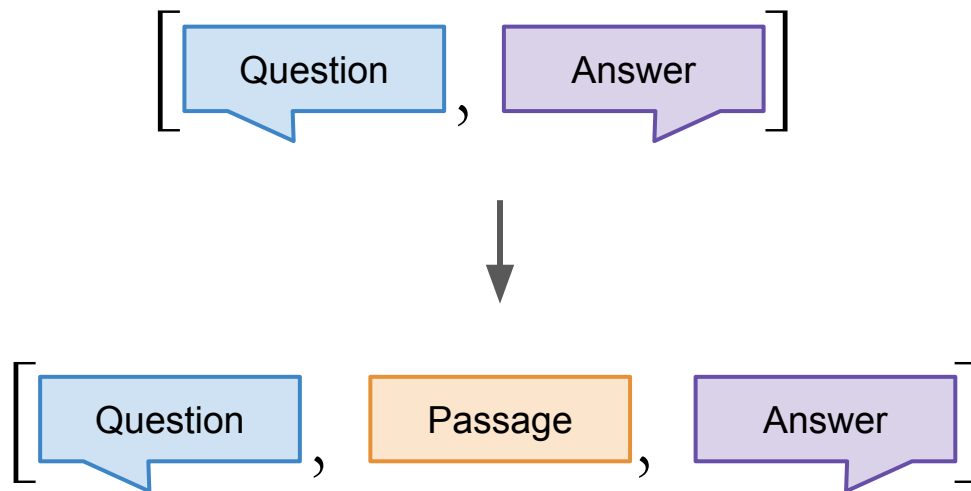
# Architecture



# Architecture



# Question-passage Dataset



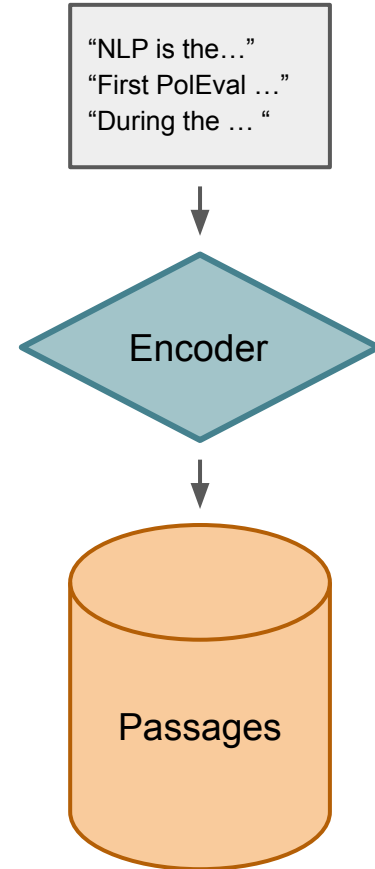
# Question-passage Dataset

- Wikipedia
  - Split into paragraphs
  - Concatenated title + paragraph
- Wiktionary
  - Proverbs
  - Idioms

“NLP is the...”  
“First PoIEval ...”  
“During the ... “

# Question-passage Dataset

- Encoders
  - Bag-of-Words
  - Universal Sentence Encoder
  - HerBERT trained on “Czy wiesz?”
  - HerBERT trained on part of data
  - HerBERT trained using answers
  - ...

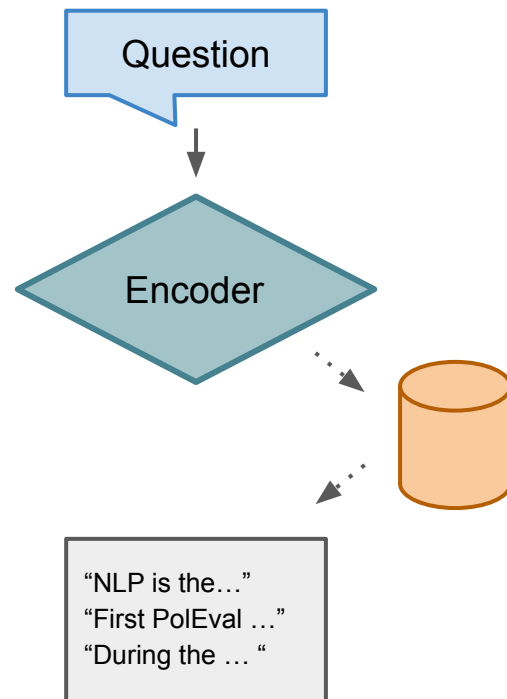




# Question-passage Dataset

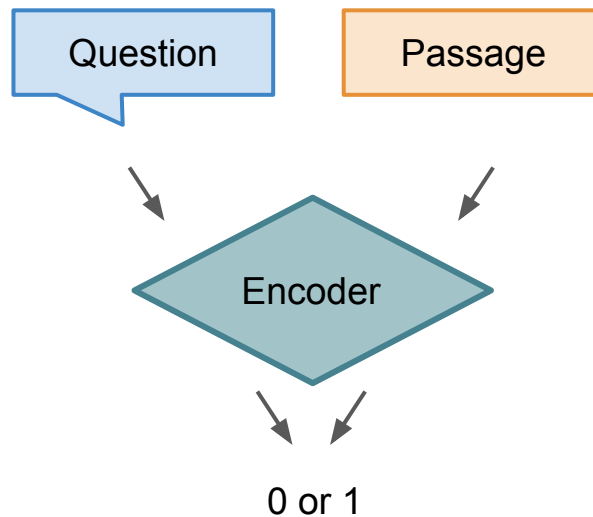
- Manually annotate candidates
  - 10k annotations
  - 2215 positives, 1347 unique questions
  - Available here:

`hf.co/datasets/allegro/polish-question-passage-pairs`

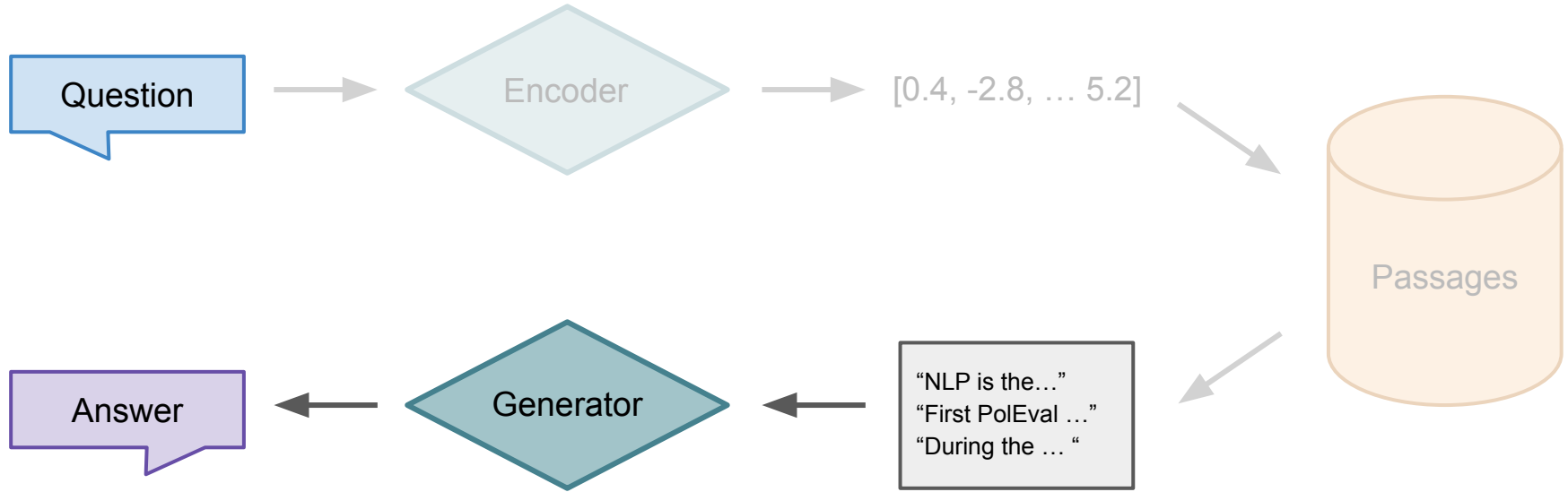


# Neural Retriever

- Architecture
  - HerBERT Large
- Dataset
  - 10k annotated pairs
  - “Czy wiesz?” dataset
- Objective
  - Contrastive loss between question & passage

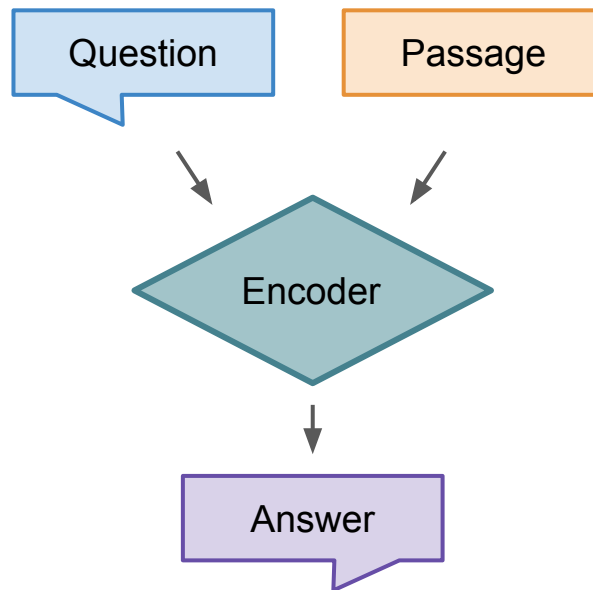


# Architecture



# Answer Generator

- Architecture
  - pIT5 Base
    - [hf.co/allegro/plt5-base](https://hf.co/allegro/plt5-base)
- First Phase
  - 1347 positives
- Second Phase
  - All question & answer pairs
  - Top 10 predicted passages



<b>Dataset slice</b>	<b>Number of questions</b>	<b>Accuracy</b>
<b>Number</b>	268	50.37
Century	54	85.19
Year	49	61.22
Other	165	35.76
<b>Binary</b>	157	66.88
Yes	97	98.97
No	60	15.00
<b>Definition</b>	417	41.97
<b>Person</b>	304	54.18
<b>Choice</b>	226	58.41
<b>Capital</b>	21	95.24
<b>Proverb</b>	10	80.00
<b>Other</b>	1124	47.06
<b>Total</b>	2500	50.96

# Thank you!

[piotr.cezary.rybak@gmail.com](mailto:piotr.cezary.rybak@gmail.com)