

# Analiza danych tekstowych

dr inż. Aleksander Smywiński-Pohl

apohllo@agh.edu.pl

<http://apohllo.pl/dydaktyka/nlp>  
(<http://apohllo.pl/dydaktyka/nlp>)

# Plan

- dane tekstowe
  - charakterystyka danych tekstowych
  - wyrażenia regularne
  - odległość edycyjna oraz inne miary odległości
  - normalizacja tekstu: tokenizacja, podział na zdania
  - wyrażenia wielosegmentowe
- podstawowe algorytmy NLP
  - stemming i lematyzacja
  - tagowanie
  - parsowanie
  - rozpoznawanie jednostek nazewniczych
  - modelowanie języka
- zastosowania

# Dane tekstowe

- charakterystyka danych tekstowych
- wyrażenia regularne
- odległość edycyjna oraz inne miary odległości
- normalizacja tekstu: tokenizacja, podział na zdania
- wyrażenia wielosegmentowe

**Ile razy w bazie danych w tabeli words w kolumnie value występuje słowo virus?**

```
select count(*) from words where value = 'virus'
```

# **Ile razy na stronie onet.pl występuje słowo "wirus"?**

- czy uwzględniamy pisownię wielką i małą literą?
- czy uwzględniamy fleksję? - wirus, wirsa, wirusowi, wirusem
- czy uwzględniamy słowotwórstwo i relacje morfologiczne? - wirsu/wirusowe, lekarz/leczenie
- jak szybko można to policzyć?

# **Ile razy w tekście występuje wyrażenie "wirusowe zapalenie żołądka"?**

- czy wszystkie słowa są odmienne?
- czy pomiędzy słowami wyrażenia mogą występować inne słowa?
- czy spacja jest zawsze tym czym nam się wydaje, że jest?

# Wyrażenia regularne - regex(p)

- pozwalają na konstrukcję wzorców dopasowujących się do tekstu
- pozwalają na uwzględnienie wariantów w pisowni
- w ograniczonym stopniu pozwalają na uwzględnienie błędów w pisowni
- w ograniczonym stopniu pozwalają na uwzględnienie fleksji
- wiele realnie działających systemów przetwarzających tekst opiera się w dużej mierze na wyrażeniach regularnych

```
In [ ]: import re

pattern = re.compile('Łódź', re.I)
text = "Łódź to stolica województwa łódzkiego"
if(pattern.search(text)):
    print("Wyrażenie zostało dopasowane")
```

```
In [ ]: tab = ["Poznań", "Łódź", "Andrychów"]
sorted(tab)
```

```
In [40]: pattern = re.compile(r'\bwirus(\w*)',)
texts = ["a _wirusów", "wirusowe zapalenie żołądka", "wirusy są groźne", "korono
wirus"]
for text in texts:
    match_data = pattern.search(text)
    if(match_data):
        print(match_data[0])
```

wirusowe  
wirusy



# Wyrażenia regularne

- Podstawowe metaznaki:
  - \* – (kwantyfikator) zero lub więcej wystąpień
  - ( ) – grupowanie wyrażeń dla kwantyfikatorów, alternatywy, dopasowań wstecznych
  - | – alternatywa – jedna opcja spośród wielu
- Kwantyfikatory (określają ile razy ma być dopasowane wyrażenie, które stoi przed nimi):
  - \* – zero lub więcej wystąpień (to samo co wyżej)
  - + – jedno lub więcej wystąpień
  - ? – zero lub jedno wystąpienie

- Kotwice (dopasowują się do pozycji w łańcuchu, a nie konkretnych znaków):
  - ^ – początek łańcucha
  - \$ – koniec łańcucha
  - \b – granica słowa
  - \< – początek słowa
  - \> – koniec słowa
- Klasy znaków:
  - [ ] – jeden ze znaków znajdujących się wewnątrz nawiasów
  - ^ – pojawiając się na początku w kontekście klasy znaków powoduje jej zanegowanie
  - a - z – zakres znaków (tylko wew. nawiasów kwadratowych)
  - \w – znak będący literą, cyfrą lub podkreśleniem
  - \s – znak będący białą spacją (spacja, tabulator, koniec linii, etc.)
  - \d – cyfra
  - . – dowolny znak (zazwyczaj – niebędący końcem linii)

<https://regex101.com/r/pTEWml/2> (<https://regex101.com/r/pTEWml/2>).

# Zaawansowane własności wyrażeń regularny

- Wsparcie dla Unicode:
  - <https://www.regular-expressions.info/unicode.html>  
(<https://www.regular-expressions.info/unicode.html>)
  - `\p{L}` - litery z dowolnego alfabetu (np. a, ą, ć, ü, 力)
  - `\p{Ll}` - mała litera z dowolnego alfabetu
  - `\p{Lu}` - wielka litera z dowolnego alfabetu
  - `\X` - dowolny znak Unicode (jak `\s`, ale automatycznie uwzględnia znaki przejścia do nowej linii)
- *positive lookahead*
  - wyrażenie `(\w+)(?= ma kota)` dopasuje się do łańcucha `Ala ma kota`, ale dopasowanie obejmie tylko słowo `Ala`.
- *negative lookbehind*
  - wyrażenie `(?!<starych ) (zł)` dopasuje się do łańcucha `10 złotych` ale nie do łańcucha `10 starych złotych`.

```
In [45]: text = "aaa"  
         re.findall("a(?=a)", text)
```

```
Out[45]: ['a', 'a']
```

# Narzędzia

- Python
  - re - standardowa biblioteka - <https://docs.python.org/3/library/re.html> (<https://docs.python.org/3/library/re.html>)
  - pip install regex - biblioteka nie wbudowana w język - <https://pypi.org/project/regex/> (<https://pypi.org/project/regex/>)
- Java
  - java.util.regex.Pattern - <https://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html> (<https://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html>)
- Ruby - <https://ruby-doc.org/core-2.5.1/Regexp.html> (<https://ruby-doc.org/core-2.5.1/Regexp.html>)
  - literał /a.\*/
- Javascript - [https://developer.mozilla.org/en-US/docs/Web/JavaScript/Guide/Regular\\_Expressions](https://developer.mozilla.org/en-US/docs/Web/JavaScript/Guide/Regular_Expressions) ([https://developer.mozilla.org/en-US/docs/Web/JavaScript/Guide/Regular\\_Expressions](https://developer.mozilla.org/en-US/docs/Web/JavaScript/Guide/Regular_Expressions))
  - literał /a.\*/

## Odległość edycyjna (odległość Levenshteina)

Minimalna liczba operacji:

- usunięcia
- dodania
- zamiany

pojedynczego znaku, dzięki którym łańcuch  $S_1$  może zostać przekształcony w łańcuch  $S_2$ .

[illegible]

źródło: [https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)  
([https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)).

		k	i	t	t	e	n
	0	1	2	3	4	5	6
s	1	<u>1</u>	2	3	4	5	6
i	2	2	<u>1</u>	2	3	4	5
t	3	3	2	<u>1</u>	2	3	4
t	4	4	3	2	<u>1</u>	2	3
i	5	5	4	3	2	<u>2</u>	3
n	6	6	5	4	3	3	<u>2</u>
g	7	7	6	5	4	4	<u>3</u>

		S	a	t	u	r	d	a	y
	0	1	2	3	4	5	6	7	8
S	1	<u>0</u>	<u>1</u>	<u>2</u>	3	4	5	6	7
u	2	1	1	2	<u>2</u>	3	4	5	6
n	3	2	2	2	3	<u>3</u>	4	5	6
d	4	3	3	3	3	4	<u>3</u>	4	5
a	5	4	3	4	4	4	4	<u>3</u>	4
y	6	5	4	4	5	5	5	4	<u>3</u>



# Implementacje

- algorytm dynamiczny - do porównywania pojedynczych łańcuchów znaków
- algorytm korekty tekstu Norviga - odwrócenie problemu <https://norvig.com/spell-correct.html> (<https://norvig.com/spell-correct.html>)
- Levenshtein automaton - znacznie szybszy niż algorytm Norviga, zaimplementowany np. w Elasticsearch

# Zastosowania miary edycyjnej

- wyszukiwanie przybliżone (fuzzy search - did you mean?)
- automatyczna korekta tekstu
- interpretacja komend głosowych

# Narzędzia

- Fuzzywuzzy - <https://github.com/seatgeek/fuzzywuzzy>  
(<https://github.com/seatgeek/fuzzywuzzy>)
- Elasticsearch - fuzzy match  
<https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-fuzzy-query.html>  
(<https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-fuzzy-query.html>)
- Algorytm Norviga - <https://norvig.com/spell-correct.html>  
(<https://norvig.com/spell-correct.html>)
- Language tool - <https://www.languagetool.org/> (<https://www.languagetool.org/>)
- biblioteka regex - <https://pypi.org/project/regex/> (<https://pypi.org/project/regex/>)

```
In [46]: pattern = re.compile('virus')
texts = ["wrius", "wirrus", "wirs"]
for text in texts:
    match_data = pattern.search(text)
    if(match_data):
        print(match_data[0])
```

```
In [48]: import regex

pattern = regex.compile(r'(virus\b){e<=1}')
texts = ["virus", "wirrus", "wirusa"]
for text in texts:
    match_data = pattern.search(text)
    if(match_data):
        print(match_data[0])
```

```
virus
wirrus
wirusa
```

# Normalizacja tekstu - tokenizacja

<https://regex101.com/r/pTEWml/3/> (<https://regex101.com/r/pTEWml/3/>).

Kwestie problematyczne:

- podział wyrażen wielosegmentowych: Bielsko-Biała
- wyrażenia z apostrofem: McDonald's
- liczby: 100 000 000, 20%
- adresu url i e-mail: <http://www.agh.edu.pl>

Narzędzia:

- ElasticSearch - <https://www.elastic.co/guide/en/elasticsearch/reference/current/analysis-tokenizers.html>  
(<https://www.elastic.co/guide/en/elasticsearch/reference/current/analysis-tokenizers.html>)
- Solr - [https://lucene.apache.org/solr/guide/6\\_6/tokenizers.html](https://lucene.apache.org/solr/guide/6_6/tokenizers.html)  
([https://lucene.apache.org/solr/guide/6\\_6/tokenizers.html](https://lucene.apache.org/solr/guide/6_6/tokenizers.html))

# Normalizacja tekstu - podział na zdania

<https://regex101.com/r/pTEWml/5/> (<https://regex101.com/r/pTEWml/5/>).

Narzędzia:

- LanguageTool, reguły SRX - <https://github.com/languagetool-org/languagetool> (<https://github.com/languagetool-org/languagetool>).
- Alternatywna implementacja dla reguł SRX - <https://github.com/loomchild/segment> (<https://github.com/loomchild/segment>).
- SpaCy <https://spacy.io/usage/linguistic-features#sbd> (<https://spacy.io/usage/linguistic-features#sbd>).

In [51]:

```
import spacy
from spacy.lang.pl import Polish

nlp = Polish()
sentencizer = nlp.create_pipe("sentencizer")
nlp.add_pipe(sentencizer)
doc = nlp("Piątek był trzecim dniem w tym tygodniu na giełdzie w Warszawie, w którym WIG20 spadł o więcej niż 4 proc. na zamknięciu sesji. W efekcie cały giełdowy tydzień zakończył się zniżką aż o ponad 15 proc., co jest trzecim najgorszym wynikiem w historii tego indeksu.")
for sent in doc.sents:
    print(sent.text)
```

Piątek był trzecim dniem w tym tygodniu na giełdzie w Warszawie, w którym WIG20 spadł o więcej niż 4 proc. na zamknięciu sesji. W efekcie cały giełdowy tydzień zakończył się zniżką aż o ponad 15 proc., co jest trzecim najgorszym wynikiem w historii tego indeksu.



# Wyrażenia wielosegmentowe

- Bielsko-Biała
- Rzeczpospolita Polska
- *post scriptum*
- Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie
- zastawka mitralna
- benzyna bezołowiowa

# Algorytmy ekstrakcji wyrażeń

- point-wise mutual information (PMI):
  - $pmi(a; b) = \log \frac{p(a,b)}{p(a)p(b)}$
- likelihood ratio-test (LLR test):
  - $llr(a; b) = 2 * (H(rows) + H(columns) - H(all))$
  - $H(cells) = - \sum_{i \in cells} (p_i * total * \log(p_i))$
  - <https://github.com/tdunning/python-llr/blob/master/llr.py>  
(<https://github.com/tdunning/python-llr/blob/master/llr.py>)
- algorytmy ML oparte o Wikipedię, np. AutoPhrase
  - <https://github.com/shangjingbo1226/AutoPhrase>  
(<https://github.com/shangjingbo1226/AutoPhrase>)
- TermoPL - oparte o miarę C, uwzględniający ograniczenia morfosyntaktyczne
  - $C - value(p)$

$$\begin{aligned} &= l(p) * \left( freq(p) - \frac{1}{r(LP)} \sum_{lp \in LP} freq(lp) \right) && r(LP) > 0 \\ & l(p) * freq(p) && r(LP) = 0 \end{aligned}$$

Forma podstawowa	Forma odmieniona	Miara "C"
prezes urząd	[prezes Urzędu]	267.73
kara pieniężny	[kara pieniężna]	69.92
zbiorowy interes konsument	[zbiorowy interes konsumentów]	54.68
ochrona konkurencja	[ochrona konkurencji]	53.33
urząd	[urząd]	30.97
przedsiębiorca	[przedsiębiorca]	30.87
mowa	[mowa]	29.60
sprawa praktyka	[sprawa praktyk]	29.00
sąd ochrona konkurencja	[sąd ochrony konkurencji]	28.00
prezes	[prezes]	27.14
postępowanie antymonopolowy	[postępowanie antymonopolowe]	25.16
droga decyzja	[droga decyzji]	23.00

Wyniki TermoPL dla Ustawy o Ochronie Konkurencji i Konsumentów

# Podstawowe algorytmy NLP

- stemming i lematyzacja
- tagowanie
- parsowanie
- rozpoznawanie jednostek nazewniczych
- modelowanie języka

# Stemming i lematyzacja

Znajdowanie formy wspólnej dla różnych form powiązanych syntaktycznie.

## Stemming (j. angielski)

This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things-names and heights and soundings-with the single exception of the red crosses and the written notes.

**Porter stemmer** (proste przekształcenia usuwające końcówki słów):

Thi wa not the map we found in Billi Bone s chest but an accur copi complet in all thing name and height and sound with the singl except of the red cross and the written note

# Lematyzacja (j. polski)

Partia polityczna jest dobrowolną organizacją, występującą pod określoną nazwą, stawiającą sobie za cel udział w życiu publicznym poprzez wywieranie metodami demokratycznymi wpływu na kształtowanie polityki państwa lub sprawowanie władzy publicznej.

## KRNNT (tagger i lematyzer)

partia polityczny być dobrowolny organizacja, występować pod określony nazwa, stawiać siebie za cel udział w życiu publiczny poprzez wywierać metoda demokratyczny wpływ na kształtować polityka państwo lub sprawować władza publiczny.

# Narzędzia

- Stempel - <http://www.getopt.org/stempel/> (<http://www.getopt.org/stempel/>),
- Morfologik - <https://github.com/morfologik/morfologik-stemming>  
(<https://github.com/morfologik/morfologik-stemming>),
- Morfeusz - <http://sgjp.pl/morfeusz/> (<http://sgjp.pl/morfeusz/>),
- KRNNT - <https://github.com/kwrobel-nlp/krnnt/> (<https://github.com/kwrobel-nlp/krnnt/>),
- Stanford NLP - [https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html)  
([https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html)),
- COMBO - <https://github.com/360er0/COMBO>  
(<https://github.com/360er0/COMBO>),

# Tagowanie morfosyntaktyczne

Partia	partia	subst:sg:nom:f
polityczna	polityczny	adj:sg:nom:f:pos
jest	być	fin:sg:ter:imperf
dobrowolną	dobrowolny	adj:sg:inst:f:pos
organizacją	organizacja	subst:sg:inst:f
,	,	interp
występującą	występować	pact:sg:inst:f:imperf:aff
pod	pod	prep:inst:nwok
określoną	określony	adj:sg:inst:f:pos
nazwą	nazwa	subst:sg:inst:f



# Tagi morfosyntaktyczne 1/2

flexeme	abbreviation	base form	example
noun	<i>subst</i>	singular nominative	<i>profesor</i>
depreciative form	<i>depr</i>	singular nominative form of the corresponding noun	<i>profesor</i>
main numeral	<i>num</i>	inanimate masculine nominative form	<i>pięć, dwa</i>
collective numeral	<i>numcol</i>	inanimate masculine nominative form of the main numeral	<i>pięć, dwa</i>
adjective	<i>adj</i>	singular nominative masculine positive form	<i>polski</i>
ad-adjectival adjective	<i>adja</i>	singular nominative masculine positive form of the adjective	<i>polski</i>
post-prepositional adjective	<i>adjp</i>	singular nominative masculine positive form of the adjective	<i>polski</i>
predicative adjective	<i>adjc</i>	singular nominative masculine positive form of the adjective	<i>zdrowy, ciekawy</i>
adverb	<i>adv</i>	positive form	<i>dobrze, bardzo</i>
non-3rd person pronoun	<i>ppron12</i>	singular nominative	<i>ja</i>
3rd-person pronoun	<i>ppron3</i>	singular nominative	<i>on</i>
pronoun SIEBIE	<i>siebie</i>	accusative	<i>siebie</i>
non-past form	<i>fin</i>	infinitive	<i>czytać</i>
future BYĆ	<i>bedzie</i>	infinitive	<i>być</i>
agglutinate BYĆ	<i>aglt</i>	infinitive	<i>być</i>
l-participle	<i>praet</i>	infinitive	<i>czytać</i>
imperative	<i>impt</i>	infinitive	<i>czytać</i>
impersonal	<i>imps</i>	infinitive	<i>czytać</i>
infinitive	<i>inf</i>	infinitive	<i>czytać</i>

<http://nkjp.pl/poligarp/help/ense2.html> (<http://nkjp.pl/poligarp/help/ense2.html>)

## Tagi morfosyntaktyczne 2/2

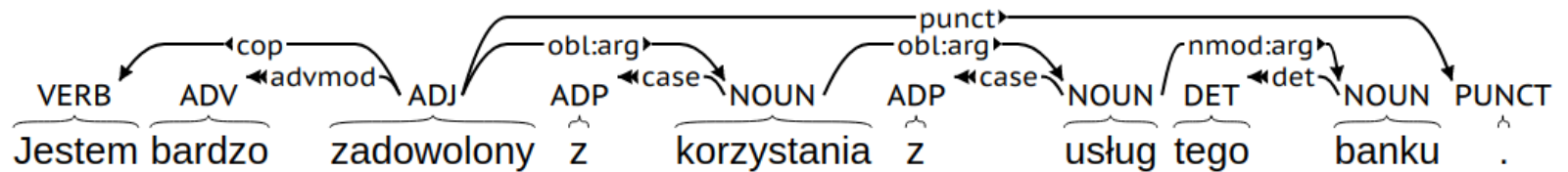
contemporary adv. participle	<i>pcon</i>	infinitive	<i>czytać</i>
anterior adv. participle	<i>pant</i>	infinitive	<i>czytać</i>
gerund	<i>ger</i>	infinitive	<i>czytać</i>
active adj. participle	<i>pact</i>	infinitive	<i>czytać</i>
passive adj. participle	<i>ppas</i>	infinitive	<i>czytać</i>
winien	<i>winien</i>	singular masculine form	<i>powinien, rad</i>
predicative	<i>pred</i>	the only form of that flexeme	<i>warto</i>
preposition	<i>prep</i>	the non-vocalic form of that flexeme	<i>na, przez, w</i>
coordinating conjunction	<i>conj</i>	the only form of that flexeme	<i>oraz</i>
subordinating conjunction	<i>comp</i>	the only form of that flexeme	<i>że</i>
particle-adverb	<i>qub</i>	the only form of that flexeme	<i>nie, -że, się</i>
abbreviation	<i>brev</i>	the full dictionary form	<i>rok, i tak dalej</i>
bound word	<i>burk</i>	the only form of that flexeme	<i>trochu, oścież</i>
interjection	<i>interj</i>	the only form of that flexeme	<i>ech, kurde</i>
punctuation	<i>interp</i>	the only form of that flexeme	<i>;, ,, (, ]</i>
alien	<i>xxx</i>	the only form of that flexeme	<i>cool, nihil</i>
unknown form	<i>ign</i>	the only form of that flexeme	

<http://nkjp.pl/poligarp/help/ense2.html> (<http://nkjp.pl/poligarp/help/ense2.html>)

# Narzędzia - taggery dla j. polskiego

- Toygger - <http://mozart.ipipan.waw.pl/~kkrasnowska/PolEval/src/>  
(<http://mozart.ipipan.waw.pl/~kkrasnowska/PolEval/src/>),
- KRNNT - <https://github.com/kwrobel-nlp/krnnt> (<https://github.com/kwrobel-nlp/krnnt>),
- COMBO - <https://github.com/360er0/COMBO>  
(<https://github.com/360er0/COMBO>),
- Stanford NLP - [https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html)  
([https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html)),
- Concraft - <https://github.com/kawu/concraft-pl>  
(<https://github.com/kawu/concraft-pl>).

# Parsing zależnościowy



# Narzędzia

- Universal dependencies - <https://universaldependencies.org/>  
(<https://universaldependencies.org/>)
- MALT Parser - <http://www.maltparser.org/> (<http://www.maltparser.org/>)
- COMBO - <https://github.com/360er0/COMBO>  
(<https://github.com/360er0/COMBO>)
- Stanford NLP - <https://stanfordnlp.github.io/stanfordnlp>  
(<https://stanfordnlp.github.io/stanfordnlp>)

```
In [52]: import stanfordnlp
#stanfordnlp.download('pl')
nlp = stanfordnlp.Pipeline(lang='pl', treebank='pl_lfg')
doc = nlp("Piątek był trzecim dniem w tym tygodniu na giełdzie w Warszawie, w którym WIG20 spadł o więcej niż 4 proc. na zamknięciu sesji.")
```

Use device: cpu

---

Loading: tokenize

With settings:

```
{'model_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg_tokenizer.pt', 'lang': 'pl', 'shorthand': 'pl_lfg', 'mode': 'predict'}
```

---

Loading: pos

With settings:

```
{'model_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg_tagger.pt', 'pretrain_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg.pretrain.pt', 'lang': 'pl', 'shorthand': 'pl_lfg', 'mode': 'predict'}
```

---

Loading: lemma

With settings:

```
{'model_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg_lemmatizer.pt', 'lang': 'pl', 'shorthand': 'pl_lfg', 'mode': 'predict'}
```

Building an attentional Seq2Seq model...

Using a Bi-LSTM encoder

Using soft attention for LSTM.

Finetune all embeddings.

[Running seq2seq lemmatizer with edit classifier]

---

Loading: depparse

With settings:

```
{'model_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg_parser.pt', 'pretrain_path': '/home/apohllo/stanfordnlp_resources/pl_lfg_models/pl_lfg.pretrain.pt', 'lang': 'pl', 'shorthand': 'pl_lfg', 'mode': 'predict'}
```

Done loading processors!

---

/pytorch/aten/src/ATen/native/LegacyDefinitions.cpp:14: UserWarning: masked\_fill\_ received a mask with dtype torch.uint8, this behavior is now deprecated, please use a mask with dtype torch.bool instead.

```
In [53]: doc.sentences[0].print_dependencies()
```

```
('Piątek', '4', 'nsubj')
('był', '4', 'cop')
('trzecim', '4', 'amod')
('dnem', '0', 'root')
('w', '7', 'case')
('tym', '7', 'det')
('tygodniu', '4', 'obl')
('na', '9', 'case')
('giełdzie', '4', 'nmod')
('w', '11', 'case')
('Warszawie', '4', 'nmod')
(',', '16', 'punct')
('w', '14', 'case')
('którym', '16', 'obl')
('WIG20', '16', 'nsubj')
('spadł', '11', 'acl:relcl')
('o', '18', 'case')
('więcej', '16', 'obl')
('niż', '21', 'case')
('4', '21', 'nummod')
('proc', '18', 'nmod')
('.', '11', 'punct')
```

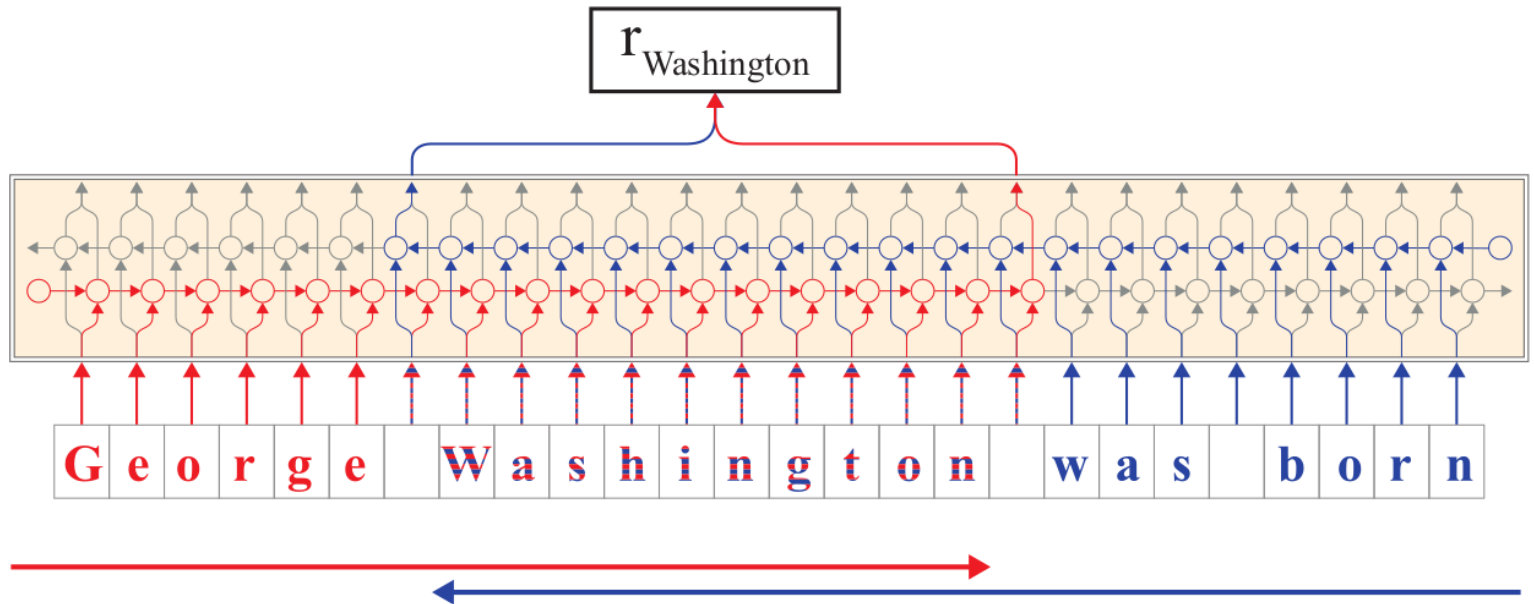


# Rozpoznawanie jednostek nazewniczych (NER)

When I told John  
PER that I wanted to move to Alaska  
LOC , he warned me that I 'd have  
trouble finding a Starbucks  
MISC there .

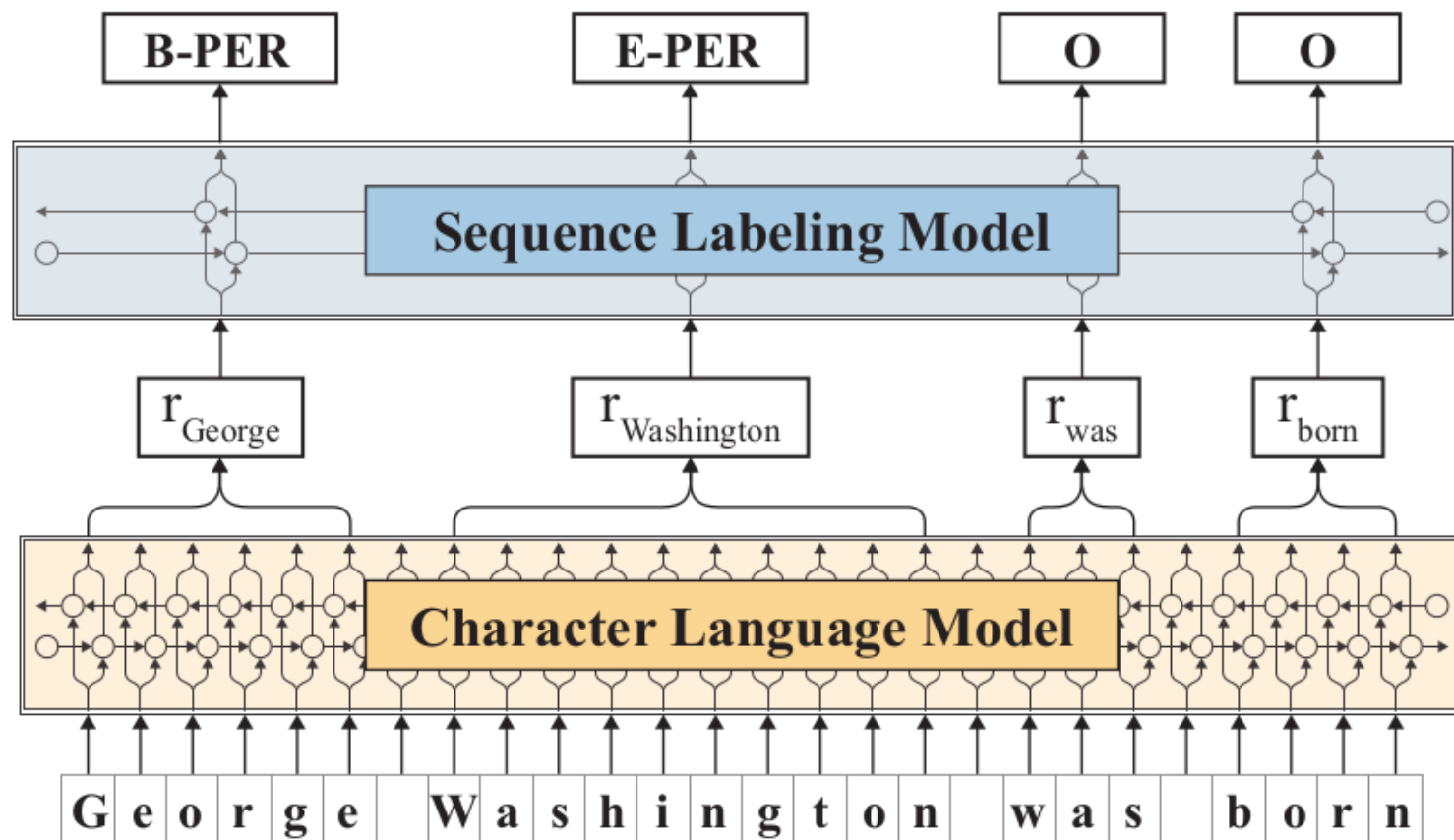
Źródło: <https://demo.allennlp.org/named-entity-recognition>  
(<https://demo.allennlp.org/named-entity-recognition>).

# Flair 1/2



Źródło: Akbik, A., Blythe, D. and Vollgraf, R., 2018. Contextual string embeddings for sequence labeling. In Proceedings of the 27th International Conference on Computational Linguistics (pp. 1638-1649).

## Flair 2/2



Źródło: Akbik, A., Blythe, D. and Vollgraf, R., 2018. Contextual string embeddings for sequence labeling. In Proceedings of the 27th International Conference on Computational Linguistics (pp. 1638-1649).



# Narzędzia

## J. polski

- Flair - <https://github.com/zalandoresearch/flair>  
(<https://github.com/zalandoresearch/flair>)
- Stanford NLP - [https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html)  
([https://stanfordnlp.github.io/stanfordnlp/installation\\_download.html](https://stanfordnlp.github.io/stanfordnlp/installation_download.html))
- PolDeepNer - <https://github.com/CLARIN-PL/PolDeepNer>  
(<https://github.com/CLARIN-PL/PolDeepNer>)
- ApplicaAI (na bazie Flair) - <https://github.com/applicaai/poleval-2018>  
(<https://github.com/applicaai/poleval-2018>)
- Liner2 - <https://github.com/CLARIN-PL/Liner2> (<https://github.com/CLARIN-PL/Liner2>)

**Narzędzia cd.**

**J. angielski**

# Modelowanie języka

	<b>i</b>	<b>want</b>	<b>to</b>	<b>eat</b>	<b>chinese</b>	<b>food</b>	<b>lunch</b>	<b>spend</b>
<b>i</b>	5	827	0	9	0	0	0	2
<b>want</b>	2	0	608	1	6	6	5	1
<b>to</b>	2	0	4	686	2	0	6	211
<b>eat</b>	0	0	2	0	16	2	42	0
<b>chinese</b>	1	0	0	0	0	82	1	0
<b>food</b>	15	0	15	0	1	4	0	0
<b>lunch</b>	2	0	0	0	0	1	0	0
<b>spend</b>	1	0	1	0	0	0	0	0

**Figure 3.1** Bigram counts for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences. Zero counts are in gray.

	<b>i</b>	<b>want</b>	<b>to</b>	<b>eat</b>	<b>chinese</b>	<b>food</b>	<b>lunch</b>	<b>spend</b>
<b>i</b>	0.002	0.33	0	0.0036	0	0	0	0.00079
<b>want</b>	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
<b>to</b>	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
<b>eat</b>	0	0	0.0027	0	0.021	0.0027	0.056	0
<b>chinese</b>	0.0063	0	0	0	0	0.52	0.0063	0
<b>food</b>	0.014	0	0.014	0	0.00092	0.0037	0	0
<b>lunch</b>	0.0059	0	0	0	0	0.0029	0	0
<b>spend</b>	0.0036	0	0.0036	0	0	0	0	0

**Figure 3.2** Bigram probabilities for eight words in the Berkeley Restaurant Project corpus of 9332 sentences. Zero probabilities are in gray.



# Generowanie tekstu

**Context (human-written):** In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**GPT-2:** The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

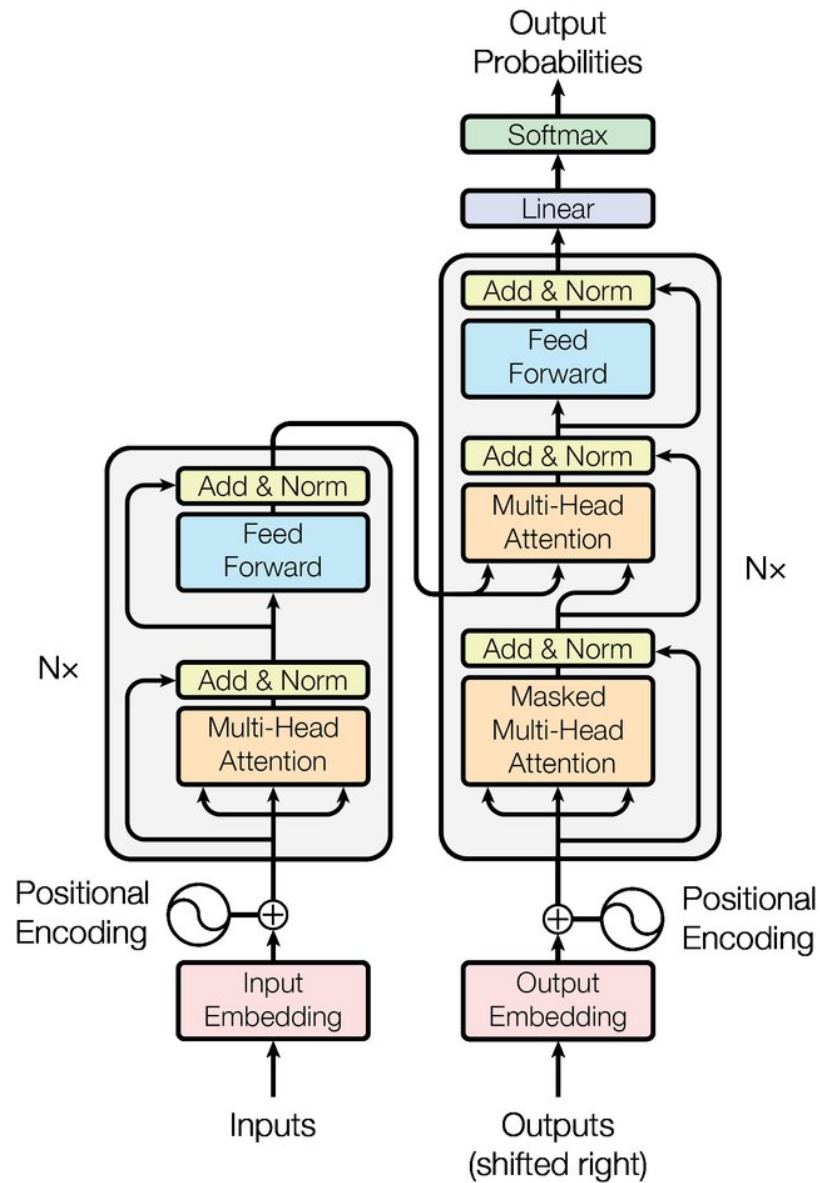
Pérez and the others then ventured further into the valley. "By the time we reached the top of one peak, the water looked blue, with some crystals on top," said Pérez.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic."

Źródło: "Language Models are Unsupervised Multitask Learners" A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever

# Architektura transformer



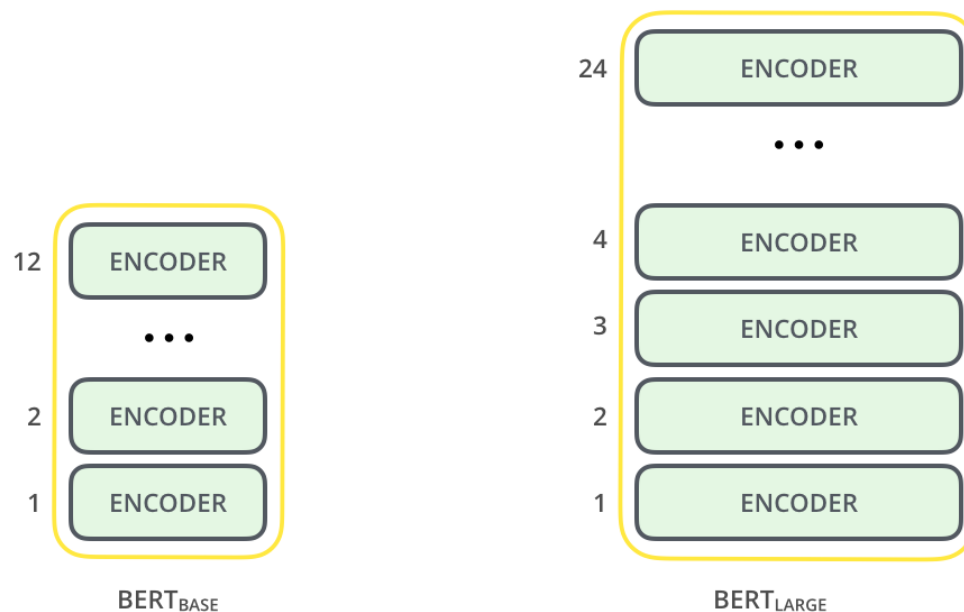
Źródło

[https://www.researchgate.net/publication/323904682\\_Tensor2Tensor\\_for\\_Neural\\_Machi](https://www.researchgate.net/publication/323904682_Tensor2Tensor_for_Neural_Machi)  
[lo=1](#)

[\(https://www.researchgate.net/publication/323904682\\_Tensor2Tensor\\_for\\_Neural\\_Mach](#)  
[lo=1\)](#)

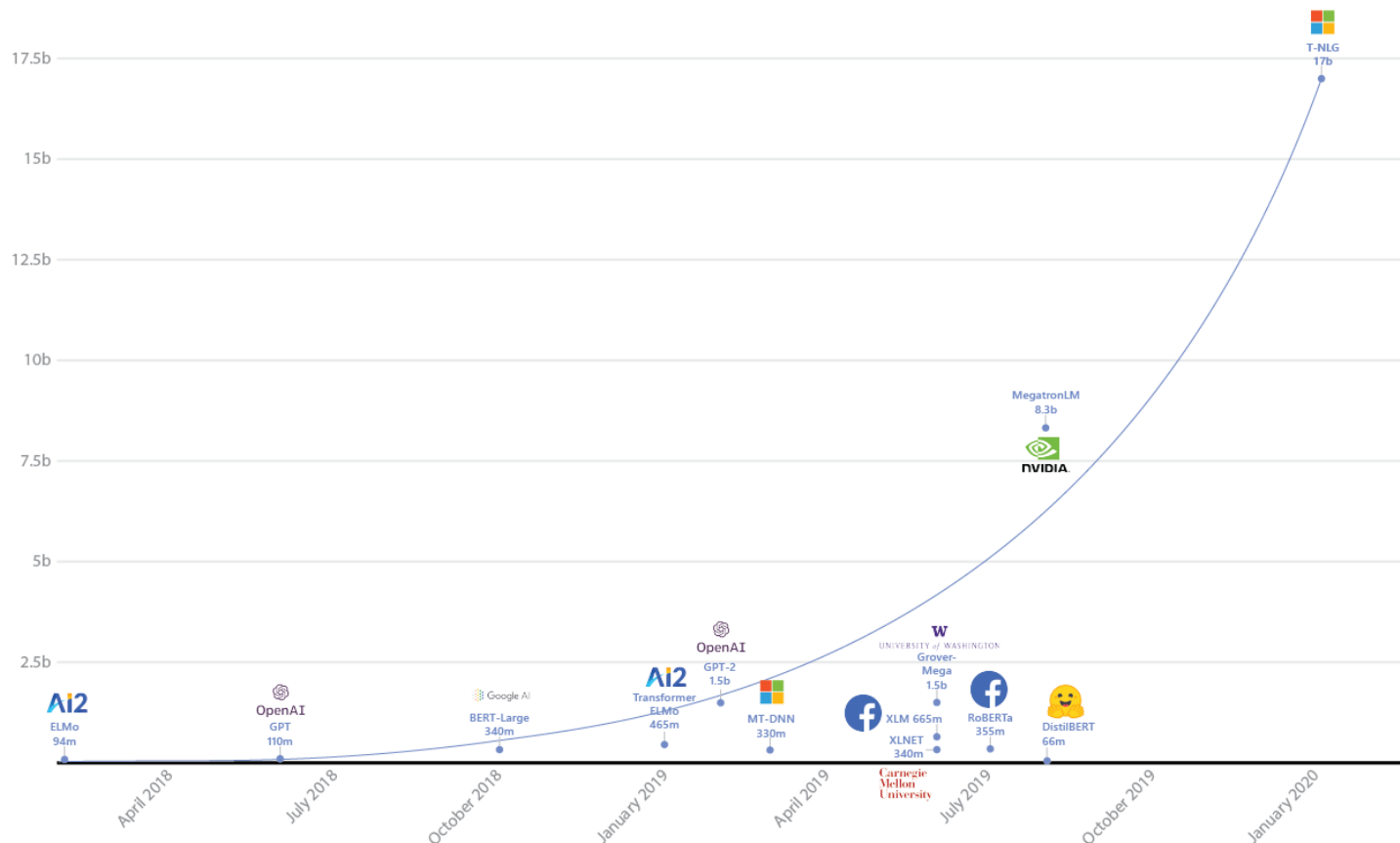
---

# Architektura BERTa



Źródło: <http://jalammar.github.io/illustrated-bert/> (<http://jalammar.github.io/illustrated-bert/>).

# Turing NLG



Źródło: <https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/> (<https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>)

## Narzędzia:

- SRI LM - <http://www.speech.sri.com/projects/srilm/>  
(<http://www.speech.sri.com/projects/srilm/>)
- Flair - <https://github.com/zalando-research/flair>  
(<https://github.com/zalando-research/flair>)
- ULMFiT - <http://nlp.fast.ai/> (<http://nlp.fast.ai/>)
- BERT - <https://github.com/google-research/bert> (<https://github.com/google-research/bert>)
- ELMo - <https://allennlp.org/elmo> (<https://allennlp.org/elmo>)

# Model T5

N=1

N=2

N=4

N=8

N=16

N=32

N=64

N=512

I love peanut butter and *bread. Thanks!! This looks delicious. I love all types of peanut butter, but especially peanut butter/jam sandwiches.*

<https://ai.googleblog.com/2020/02/exploring-transfer-learning-with-t5.html>  
(<https://ai.googleblog.com/2020/02/exploring-transfer-learning-with-t5.html>)



**Zastosowania**

# Praca z danymi tekstowymi

1. zebranie danych tekstowych
2. normalizacja danych tekstowych
3. anotacja danych
4. budowa modelu
5. testy
6. wdrożenie

# Papers with code

# Natural Language Processing

📈 429 leaderboards • 232 tasks • 100 datasets • 3560 papers with code

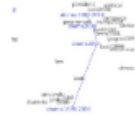
## Representation Learning



### Representation Learning

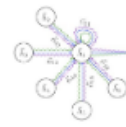
📈 3 leaderboards

439 papers with code



### Word Embeddings

426 papers with code



### Graph Embedding

90 papers with code

► [See all 17 tasks](#)

## Machine Translation

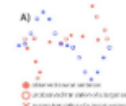


### Machine Translation

📈 44 leaderboards



### Transliteration



### Unsupervised Machine Translation

📈 9 leaderboards

# SQuAD

## Steam\_engine

### The Stanford Question Answering Dataset

Steam engines are external combustion engines, where the working fluid is separate from the combustion products. Non-combustion heat sources such as solar power, nuclear power or geothermal energy may be used. The ideal thermodynamic cycle used to analyze this process is called the Rankine cycle. In the cycle, water is heated and transforms into steam within a boiler operating at a high pressure. When expanded through pistons or turbines, mechanical work is done. The reduced-pressure steam is then condensed and pumped back into the boiler.

**Along with geothermal and nuclear, what is a notable non-combustion heat source?**

*Ground Truth Answers:* solar solar power solar power, nuclear power or geothermal energy solar

**What ideal thermodynamic cycle analyzes the process by which steam engines work?**

*Ground Truth Answers:* Rankine Rankine cycle Rankine cycle Rankine cycle

źródło: [https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/Steam\\_engine.html](https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/Steam_engine.html) ([https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/Steam\\_engine.html](https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/Steam_engine.html)).

## Zasoby