

## Przetwarzanie Języków Naturalnych

projekt pod tytułem

# Generacja tytułów publikacji naukowych

W ostatnich czasach w środowisku naukowym pojawia się coraz więcej prac naukowych tematycznie związanych z uczeniem maszynowym, sieciami neuronowymi i szeroko pojętą sztuczną inteligencją. Możliwość przewidywania tempa rozwoju badań jest bardzo istotną umiejętnością. Z tego powodu wykorzystać można dostępne obecnie narzędzia do przewidywania tego jakie będą prace naukowe które powstaną w przyszłości, a w rozważanym przypadku tytuły tych prac.

W celu przybliżenia procesu jaki należy przeprowadzić, by przewidzieć tytuły prac naukowych posłużę się przykładem programu dostępnego pod linkiem <https://github.com/csinva/gpt-paper-title-generator> którego autorami są Chandan Singh oraz Jamie Simon. Skonstruowany przez nich program miał za zadanie, wykorzystując ogólnodostępne źródła prac naukowych oraz narzędzia takie jak modele generatywne aby, wygenerować tytuły prac naukowych jakie mogą pojawić się w przyszłości.

Pierwszym etapem jest pozyskanie danych w postaci tytułów artykułów naukowych. Bazą danych która została wykorzystana jest witryna arXiv na której znajdują się około 2 250 000 artykułów naukowych z takich dziedzin jak fizyka, matematyka, informatyka i inne pokrewne nauki.



arXiv is a free distribution service and an open-access archive for 2,254,199 scholarly articles in the fields of physics, mathematics, computer science, quantitative biology, quantitative finance, statistics, electrical engineering and systems science, and economics. Materials on this site are not peer-reviewed by arXiv.

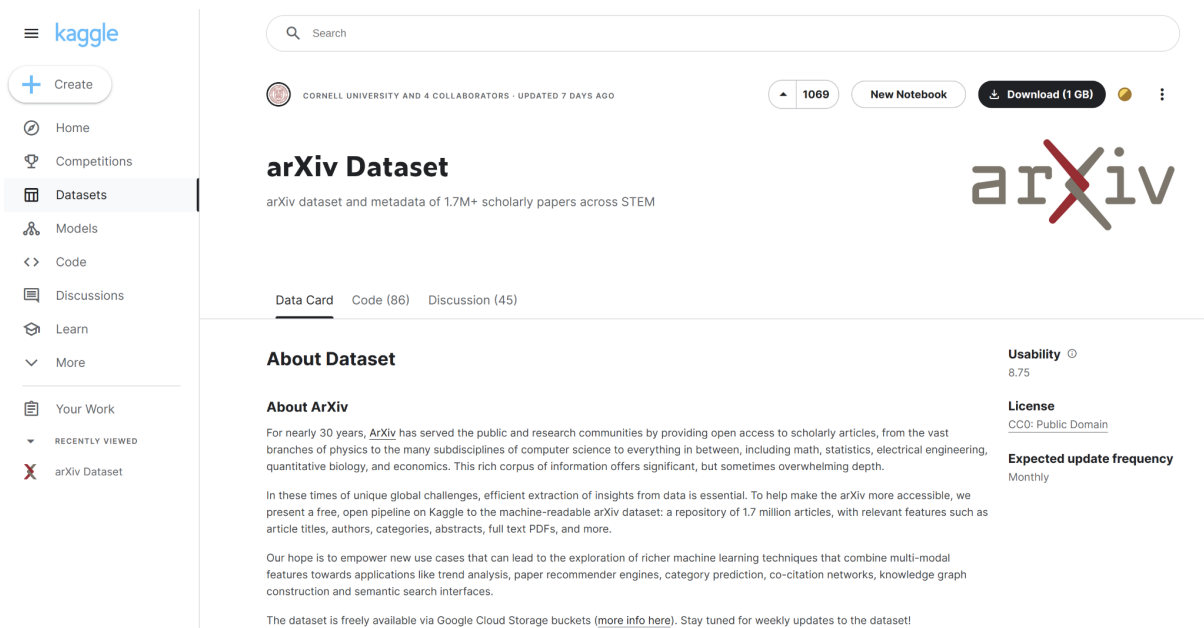
Subject search and browse:

Physics

### News

Read about recent news and updates on [arXiv's blog](#). (View the former "what's new" pages here).

Dane potrzebne do analizy zostały pobrane z witryna kaggle która udostępnia darmowe datasey które można wykorzystać we własnych programach. W tym przypadku pobraliśmy dataset składający się z artykułów dostępnych na arXiv



The screenshot shows the Kaggle interface for the 'arXiv Dataset'. The left sidebar contains navigation options like 'Home', 'Competitions', 'Datasets', 'Models', 'Code', 'Discussions', 'Learn', 'More', 'Your Work', and 'RECENTLY VIEWED'. The main content area displays the dataset title 'arXiv Dataset' with a subtitle 'arXiv dataset and metadata of 1.7M+ scholarly papers across STEM'. Below this, there are tabs for 'Data Card', 'Code (86)', and 'Discussion (45)'. The 'About Dataset' section provides a detailed description of the dataset's origin and purpose. On the right, there are metrics for 'Usability' (8.75), 'License' (CC0: Public Domain), and 'Expected update frequency' (Monthly).

Dane te pobrać można bezpośrednio ze strony lub wykorzystać aplikację konsolową kaggle

```
! pip install kaggle
! kaggle datasets download Cornell-University/arxiv
! unzip arxiv.zip
✓ 1m 52.2s

Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: kaggle in /home/szymon/.local/lib/python3.8/site-packages (1.5.13)
Requirement already satisfied: certifi in /usr/lib/python3/dist-packages (from kaggle) (2019.11.28)
Requirement already satisfied: python-dateutil in /home/szymon/.local/lib/python3.8/site-packages (from kaggle) (2.8.2)
Requirement already satisfied: urllib3 in /usr/lib/python3/dist-packages (from kaggle) (1.25.8)
Requirement already satisfied: python-slugify in /home/szymon/.local/lib/python3.8/site-packages (from kaggle) (8.0.1)
Requirement already satisfied: six>=1.10 in /usr/lib/python3/dist-packages (from kaggle) (1.14.0)
Requirement already satisfied: tqdm in /home/szymon/.local/lib/python3.8/site-packages (from kaggle) (4.65.0)
Requirement already satisfied: requests in /home/szymon/.local/lib/python3.8/site-packages (from kaggle) (2.28.2)
Requirement already satisfied: text-unidecode>=1.3 in /home/szymon/.local/lib/python3.8/site-packages (from python-slugify->kaggle) (1.3)
Requirement already satisfied: idna<4, >=2.5 in /usr/lib/python3/dist-packages (from requests->kaggle) (2.8)
Requirement already satisfied: charset-normalizer<4, >=2 in /home/szymon/.local/lib/python3.8/site-packages (from requests->kaggle) (3.1.0)

[notice] A new release of pip is available: 23.0.1 -> 23.1.2
[notice] To update, run: python3 -m pip install --upgrade pip
Downloading arxiv.zip to /home/szymon/PJN/projekt
100%|#####| 1.15G/1.15G [01:14<00:00, 13.6MB/s]
100%|#####| 1.15G/1.15G [01:15<00:00, 16.4MB/s]
Archive: arxiv.zip
  inflating: arxiv-metadata-oai-snapshot.json
```

Archiwum waży około 1.2 GB, natomiast po rozpakowaniu rozmiar wynosi około 3.5GB zawierając ponad 2 200 000 rekordów w postaci JSON.

```
> ls -lh
total 4.7G
-rw-rw-r-- 1 szymon szymon 52K May 13 14:14 1.ipynb
-rw-rw-r-- 1 szymon szymon 3.5G May 6 23:52 arxiv-metadata-oai-snapshot.json
-rw-rw-r-- 1 szymon szymon 1.2G May 13 14:37 arxiv.zip
~/PJN/projekt >
```

```

import json

attributes = ["authors", "title", "update_date", "categories"]

papers = []
count = 0
with open('arxiv-metadata-oai-snapshot.json') as f:
    for line in f:
        count += 1
        if count % 100000 == 0:
            print(f'on paper {count}')
        json_object = json.loads(line)
        paper = {}
        for attr in attributes:
            paper[attr] = json_object[attr]
        papers.append(paper)

import pandas as pd
df = pd.DataFrame(papers)

df
✓ 32.8s

```

Pobrane dane w postaci JSON zostały przetworzone na ramkę danych biblioteki pandas w celu dalszej obróbki

	authors	title	update_date	categories
0	C. Ball'vazs, E. L. Berger, P. M. Nadolsky, C.-...	Calculation of prompt diphoton production cros...	2008-11-26	hep-ph
1	Ileana Streinu and Louis Theran	Sparsity-certifying Graph Decompositions	2008-12-13	math.CO cs.CG
2	Hongjun Pan	The evolution of the Earth-Moon system based o...	2008-01-13	physics.gen-ph
3	David Callan	A determinant of Stirling cycle numbers counts...	2007-05-23	math.CO
4	Wael Abu-Shammala and Alberto Torchinsky	From dyadic $\lambda_{\alpha}$ to $\lambda_{\alpha}$	2013-10-15	math.CA math.FA
...	...	...	...	...
2250218	R. Prozorov, M. Konczykowski, B. Schmidt, Y. Y...	On the origin of the irreversibility line in t...	2009-10-30	supr-con cond-mat.supr-con
2250219	Durga P. Choudhury, Balam A. Willemsen, John S...	Nonlinear Response of HTSC Thin Film Microwave...	2016-11-18	supr-con cond-mat.supr-con
2250220	Balam A. Willemsen, J. S. Derov and S.Sridhar ...	Critical State Flux Penetration and Linear Mic...	2009-10-30	supr-con cond-mat.supr-con
2250221	Yasumasa Hasegawa (Himeji Institute of Technol...	Density of States and NMR Relaxation Rate in A...	2009-10-30	supr-con cond-mat.supr-con
2250222	Naoki Enomoto, Masanori Ichioka and Kazushige ...	Ginzburg Landau theory for d-wave pairing and ...	2009-10-30	supr-con cond-mat.supr-con

2250223 rows x 4 columns

Jak widzimy, ramka danych zawiera około 2 500 000 rekordów prac naukowych ze wszystkich dziedzin dostępnych na platformie arXiv

```

ai_papers = [paper for paper in papers if "cs.AI" in paper["categories"] or "stat.ML" in paper["categories"] or "cs.LG" in paper["categories"]]
df = pd.DataFrame(ai_papers)

import re
regex = re.compile(" and |,")

df["authors"] = [[str(a).strip() for a in regex.split(author) if len(str(a).strip()) > 3] for author in df["authors"]]
df = df.sort_values(by="update_date", ascending=False)

df
✓ 3.2s

```

	authors	title	update_date	categories
168839	[Shuting Shen, Xi Chen, Ethan X. Fang, Junwei Lu]	Combinatorial Inference on the Optimal Assortm...	2023-05-05	stat.ML cs.LG
179649	[Ming-Kun Xie, Jia-Hao Xiao, Gang Niu, Masashi...	Class-Distribution-Aware Pseudo Labeling for S...	2023-05-05	cs.LG
169390	[Alon Albalak, Colin Raffel, William Yang Wang]	Improving Few-Shot Generalization by Exploring...	2023-05-05	cs.LG cs.CL
179618	[UngJin Na, Moonhee Choi, HangJin Jo]	Critical heat flux diagnosis using conditional...	2023-05-05	physics.flu-dyn cs.LG
171498	[Soham Rohit Chitnis, Sidong Liu, Tirtharaj Da...	Domain-Specific Pre-training Improves Confiden...	2023-05-05	cs.CV cs.AI cs.LG
...	...	...	...	...
180173	[Petro M. Gopych]	A Neural Network Assembly Memory Model Based o...	2007-05-23	cs.AI cs.IR cs.NE q-bio.NC q-bio.QM
180172	[Peter D. Turney, Michael L. Littman, Jeffrey ...	Combining Independent Modules to Solve Multipl...	2007-05-23	cs.CL cs.IR cs.LG
180171	[Peter D. Turney (National Research Council of...	Measuring Praise and Criticism: Inference of S...	2007-05-23	cs.CL cs.IR cs.LG
180170	[Wolfgang Mayer, Markus Stumptner]	Model-Based Debugging using Multiple Abstract ...	2007-05-23	cs.SE cs.AI
1	[T. Kosel, I. Grabec]	Intelligent location of simultaneously active ...	2007-05-23	cs.NE cs.AI

181083 rows x 4 columns

Następnie dane te zostały odfiltrowane tak, aby pozostały jedynie rekordy których tematyką jest sztuczna inteligencja czy uczenie maszynowe. Po filtracji zostało około 180 tys. rekordów.

```

authors_with_papers = {}

for paper in df.values:
    authors = paper[0]
    title = paper[1]
    for author in authors:
        if author in authors_with_papers:
            authors_with_papers[author].append(title)
        else:
            authors_with_papers[author] = [title]

output = pd.DataFrame([{"author": author, "papers": authors_with_papers[author], "paper_count": len(authors_with_papers[author])} for author in
output.sort_values(by="paper_count", ascending=False)

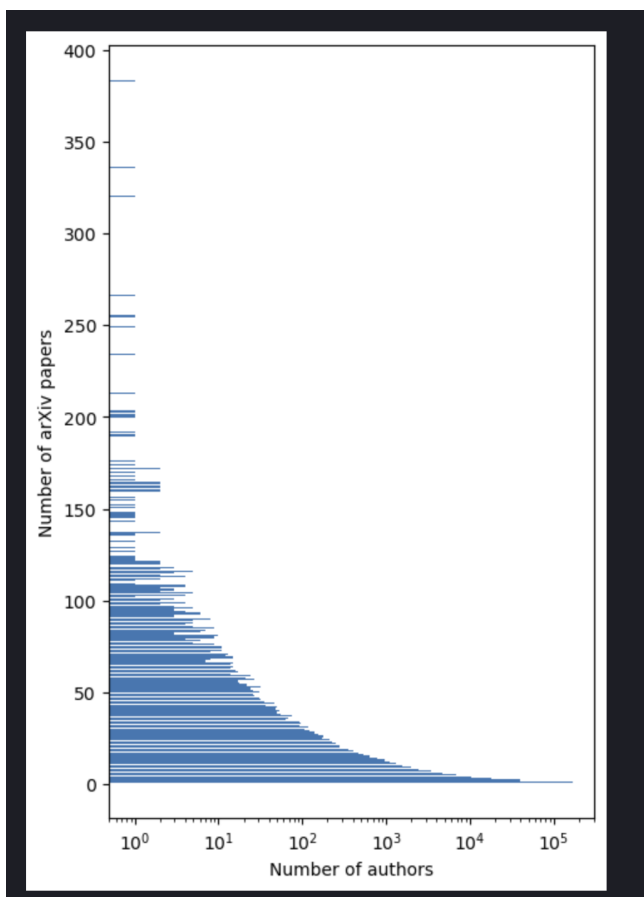
```

✓ 2.9s Python

	author	papers	paper_count
10152	Yoshua Bengio	[Hyena Hierarchy: Towards Larger Convolutional...	383
2130	Yang Liu	[Performative Prediction with Bandit Feedback...	336
3403	Sergey Levine	[Offline RL for Natural Language Generation wi...	320
4749	Bo Li	[Model Explainability in Physiological and Hea...	266
15088	Michael I. Jordan	[Online Learning in Stackelberg Games with an ...	255
...	...	...	...
128610	Igor Khokhlov	[Attacks, Defenses, And Tools: A Framework To ...	1
128611	Mehdi Mirakhorli	[Attacks, Defenses, And Tools: A Framework To ...	1
128616	Songqing Yue	[Imbalanced Malware Images Classification: a C...	1
128619	Zhengren Wang	[Listing Maximal k-Plexes in Large Real-World ...	1
267544	Markus Stumptner	[Model-Based Debugging using Multiple Abstract...	1

267545 rows × 3 columns

Po zliczeniu ilości prac według konkretnych pojedynczych autorów otrzymano następujące wyniki



Jak widać na wykresie, stosunek ilości prac ilości autorów którzy napisali tę pracę ma charakter wykładniczy.

Po przygotowaniu danych, zostały one wprowadzone do modelu GPT w formie pięciu kolejnych prac konkretnego autora w oczekiwaniu na to, że model przewidzi jakie tematy mogą mieć kolejne pięć prac tego autora.

## 1 Author-specific paper titles (prompting gpt3)

To generate author-specific titles, we take the five most recent titles from each author with atleast 3 arXiv AI papers (cs.ML, cs.LG, stat.ML). We then format the papers using the following template and query for a new title using GPT-3:

Here is a list of related machine-learning papers:

```
> [title 1]
> [title 2]
...
> [title 5]
> _____
```

```
for i, author in enumerate(tqdm(authors)):
    if not author in authors_save:
        query = prompt + '\n> '.join(authors_dict_titles[author][-papers_in_context:]) + '\n>'
        completion = openai.Completion.create(
            engine="text-davinci-002", prompt=query,
            n=gens_per_author, stop='>'
        )
        authors_save[author] = [completion.choices[i].text for i in range(len(completion.choices))]
    if i % 200 == 0:
        pickle.dump(authors_save, open(f'gen_titles/authors_save_{i}.pkl', 'wb'))
pickle.dump(authors_save, open(f'gen_titles/authors_save_full.pkl', 'wb'))
```

Chandan Singh, 7 months

Wykorzystany został model GPT-3 od OpenAI, a do komunikacji z modelem posłużyło oficjalne Pythonowe API które pozwala na prostą komunikację z modelem. Konkretnie został wykorzystany model “text-davinci-002”, w trybie uzupełniania wprowadzonego tekstu.

Wprowadzając pięć następujących tytułów

```
> Hierarchical Shrinkage: improving the accuracy and interpretability of tree-based methods
> Fast Interpretable Greedy-Tree Sums (FIGS)
> Adaptive wavelet distillation from neural networks through interpretations
> Emb-GAM: an Interpretable and Efficient Predictor using Pre-trained Language Models
> Explaining Patterns in Data with Language Models via Interpretable Autoprompting
> _____
```

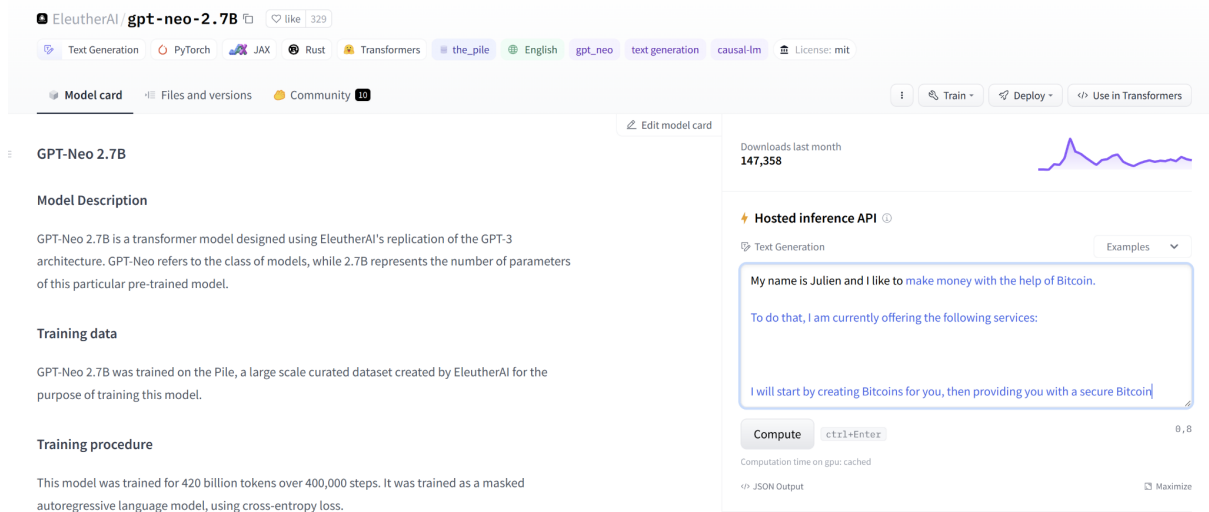
Otrzymano następujący wynik

1. Towards Interpretable Natural Language Processing: A Survey
2. A Unified Framework for Interpretable Machine Learning
3. Compositional Attention Networks for Machine Reasoning
4. Achieving Open Vocabulary Neural Machine Translation
5. A Deep Understanding of Neural Networks through Deep Visualization

Zauważono jednak, że otrzymane wyniki nie mają bezpośredniego związku z podawanymi tytułami, a jedynie z popularnymi w ostatnim czasie tematami prac naukowych.

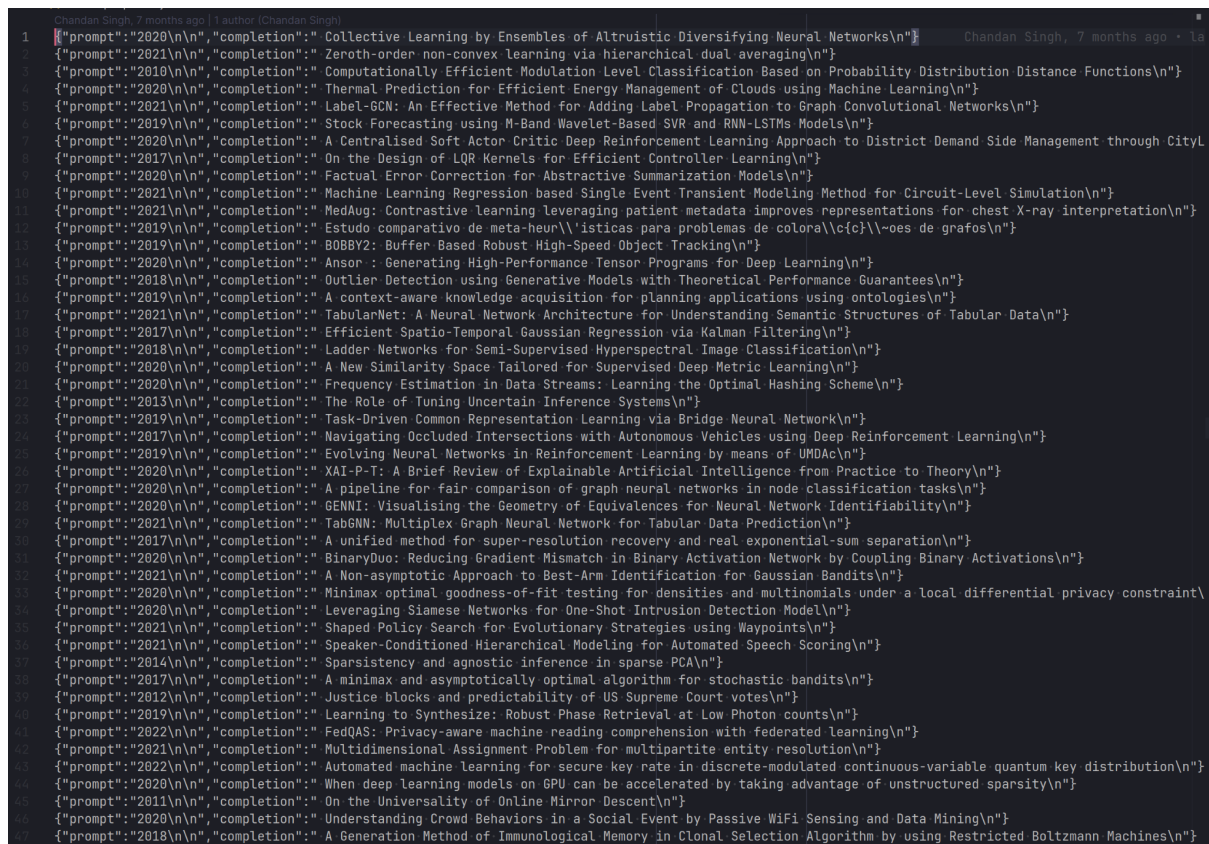
Z tego powodu postanowiono nauczyć model konkretnymi przykładami tematów prac, które pojawiały się po okresie uczenia się modelu.

Jako model bazowy wykorzystano model gpt-neo z 2.7 mld parametrów.



The screenshot shows the Hugging Face model card for GPT-Neo 2.7B. The card includes a 'Model Description' section stating it's a transformer model designed using EleutherAI's replication of the GPT-3 architecture. The 'Training data' section mentions it was trained on the Pile dataset. The 'Training procedure' section notes it was trained for 420 billion tokens over 400,000 steps. On the right, there's a 'Hosted inference API' section with a text input field containing a prompt: 'My name is Julien and I like to make money with the help of Bitcoin. To do that, I am currently offering the following services: I will start by creating Bitcoins for you, then providing you with a secure Bitcoin'. Below the input is a 'Compute' button and a 'JSON Output' checkbox.

Został on nauczony danymi spreparowanymi w następujący sposób. Jak widzimy model potrzebuje parametru prompt który oznacza tekst wpisywany przez użytkownika w momencie kiedy chce uzyskać autouzupełnianie, natomiast completion oznacza odpowiedź którą chcielibyśmy uzyskać



The screenshot shows a code editor with a list of prompt-completion pairs. Each pair is a JSON object with 'prompt' and 'completion' keys. The prompts are various research titles and technical topics, and the completions are their corresponding titles. For example, the first pair is: {"prompt": "2020", "completion": "Collective Learning by Ensembles of Altruistic Diversifying Neural Networks"}. The list continues with many other titles, such as "Zeroth-order non-convex learning via hierarchical dual averaging", "Computationally Efficient Modulation Level Classification Based on Probability Distribution Distance Functions", and "A pipeline for fair comparison of graph neural networks in node classification tasks".

Uzyskany model został umieszczony w sieci i jest ogólnie dostępny

csinva / **gpt-neo-2.7B-titles** like 0

Text Generation PyTorch Transformers gpt\_neo License: apache-2.0

Model card Files and versions Community

Full code and details at <https://github.com/csinva/gpt-paper-title-generator>

**Model**

- finetunes starting from the [gpt-neo-2.7B checkpoint](#)
  - for training details see [the training script](#)
- inference
  - should prepend with a year and two newlines before querying for a title, e.g. 2022\n\n

```
from transformers import AutoModelForCausalLM, pipeline, AutoTokenizer
model = AutoModelForCausalLM.from_pretrained("csinva/gpt-neo-2.7B-titles")
tokenizer = AutoTokenizer.from_pretrained("EleutherAI/gpt-neo-2.7B")
pipe = pipeline('text-generation', model=model, tokenizer=tokenizer)
pipe('2022\n\n')
```

Downloads last month: 2

**Hosted inference API**

Text Generation Examples

2021

Compute ctrl+Enter 0,2

This model can be loaded on the Inference API on-demand.

JSON Output Maximize

Następnie model poproszono o wygenerowanie tytułów prac naukowych dla poszczególnych lat. Wyniki prezentują się następująco:

### 2022

- Diverse Datasets for Learning to Rank via Knowledge Graph Embedding
- Machine learning-driven method for high-throughput single-cell analysis of differentiation and lineage commitment
- On the Sample Complexity of Differentially Private Learning
- Data-Dependent Weight Normalization for Improved Image Resolution
- Adaptive Densely Connected Networks for the Generation and Visualization of Object Deformations
- Exploring the Implicit Bias in Transfer Learning using Imitation Learning

### 2023 (These samples tend to just be similar to 2021/2022 where the majority of the training data lies)

- An Interpretable Dynamic Network for Spatiotemporal Pattern Prediction in High-Dimensional Time Series Data
- Multimodal Deep Learning for Automated Cancer Histopathology Analysis
- Reinforced Learning for Robust and Accurate Object Detection
- A Machine Learning Approach to High Sensitivity Data Processing
- Adversarial Robustness for Graph Neural Networks in Network Intrusion Detection
- Reinforcement Learning via Exploration and Rehearsal for Learning from Demonstrations

### 2010 (Seems to properly generate older titles)

- Learning in a Dynamic, Clustered and Homogeneous Heterogeneous Markov Decision Process
- An Empirical Analysis of the Regularization of the Gaussian Process Regression
- A Scalable Clustering Algorithm under Heterogeneous Data
- A Unified Representation for Probabilistic Time Series Forecasting
- A Hybrid Approach to Automatic Alignment and Localization of Digital Object Platforms
- Bayesian nonparametric modeling of random fields

Wygenerowane tytuły z pewnością przypominają rzeczywiste tytuły prac naukowych.

Postanowiono sprawdzić dokładność wygenerowanych tytułów i porównać je z rzeczywistymi tytułami powstałymi w tym okresie. Do tego celu postanowiono wykorzystać algorytm BLEU który powstał do porównywania automatycznego tłumaczenia tekstu z tekstem przetłumaczonym przez profesjonalnego tłumacza.

## BLEU [\[edytuj\]](#)

🌐 10 języków ▾

Artykuł [Dyskusja](#)

[Czytaj](#) [Edytuj](#) [Edytuj kod źródłowy](#) [Wyświetl historię](#) [Narzędzia](#) ▾

**BLEU (Bilingual Evaluation Understudy)** – algorytm do ewaluacji jakości tłumaczenia automatycznego z jednego języka naturalnego na inny. Jakość jest rozumiana jako korelacja między danymi wyjściowymi a tekstem ludzkim: „im bliższe tłumaczenie automatyczne jest profesjonalnemu tłumaczeniu ludzkiemu, tym jest lepsze”<sup>[1]</sup>. BLEU był jedną z pierwszych metryk, która uzyskała wysoką korelację z ludzkim osądem jakości<sup>[1][2]</sup>. Pozostaje także najbardziej popularną z metod.

Punkty liczone są dla pojedynczych przetłumaczonych segmentów – zwykle zdań – przez porównanie ich ze zbiorem tłumaczeń referencyjnych dobrej jakości. Punkty te są następnie uśredniane w obrębie całego korpusu, aby oszacować całkowitą jakość tłumaczenia. Pod uwagę nie są brane zrozumiałość oraz poprawność gramatyczna.

BLEU jest zaprojektowany, aby przybliżać ludzką ocenę na poziomie dużych korpusów i nie sprawdza się do oceny pojedynczych zdań.

Ewaluując wygenerowane tytuły postarano się znaleźć tytuły im odpowiadające. Oto pięć pierwszych wyników:

A	B
Understanding the effect of data augmentation in generative adversarial networks	Understanding the effect of data augmentation in self-supervised anomaly detection
Adversarial attacks on graph neural networks	Sparse vicious attacks on graph neural networks
Differentiable reinforcement learning for continuous control	Normality-guided distributional reinforcement learning for continuous control
Multilevel representation learning for time series forecasting	Out-of-distribution representation learning for time series classification
Unsupervised feature learning for medical image segmentation	Distributed contrastive learning for medical image segmentation

Wygenerowane tytuły są praktycznie nie do odróżnienia od tych które stanowią nazwy rzeczywistych prac naukowych. Autorzy zwracają jednak uwagę na fakt, że wygenerowane treści są bardziej ogólne od odpowiadającym im rzeczywistym tytułom.

Przedstawiona w tym raporcie analiza pokazuje potencjał drzemiący w językowych modelach generatywnych w kontekście przewidywania tudzież generowania treści na podstawie danych historycznych, co może być przydatne przy próbach ustalenia kierunku rozwoju różnych technologii czy gałęzi nauki. Już teraz istnieją narzędzia które połączone ze sobą dają zaskakująco świetne rezultaty.