# Project Intermediate Report 2

> Project GitHub link

## Team Members

- Rafail Venediktov

- Aliaksei Korshuk

- Panov Evgenii

## Current Progress

### Dataset

To test our hypotheses, we decided to take a shortened version of the huge dataset as **Fraser/short-jokes** dataset

The **Fraser/short-jokes** dataset available on the Hugging Face website is a collection of short, humorous jokes in English language. This dataset contains over 200,000 jokes, which are categorized into different categories such as puns, one-liners, knock-knock jokes, and more.

Each joke in the dataset is represented as a text string and is associated with its corresponding category. The dataset can be used to train and test machine learning models for various natural language processing tasks such as joke generation, humor detection, and sentiment analysis.

This dataset is particularly useful for developing models that can generate humorous content, as it provides a large amount of diverse and high-quality examples of short jokes. Additionally, the dataset can be used for other applications such as building chatbots or language-based recommendation systems that incorporate humor into their responses.

### Model

We decided to use the **<u>GPT-2</u>** model

Pre-trained language model based on the GPT-2 architecture, fine-tuned on a large dataset of jokes. This model is capable of generating humorous and entertaining jokes in response to input prompts or queries.

The GPT-2 architecture is a state-of-the-art deep learning model that is known for its ability to generate high-quality, coherent text that closely mimics human language. The fine-tuning of the model on a dataset of jokes makes it particularly well-suited for generating humor-related content.

The model can be used in a wide range of applications, such as chatbots, virtual assistants, and recommendation systems that require generating humorous content. Additionally, it can be used as a starting point for further fine-tuning on specific domains or datasets.
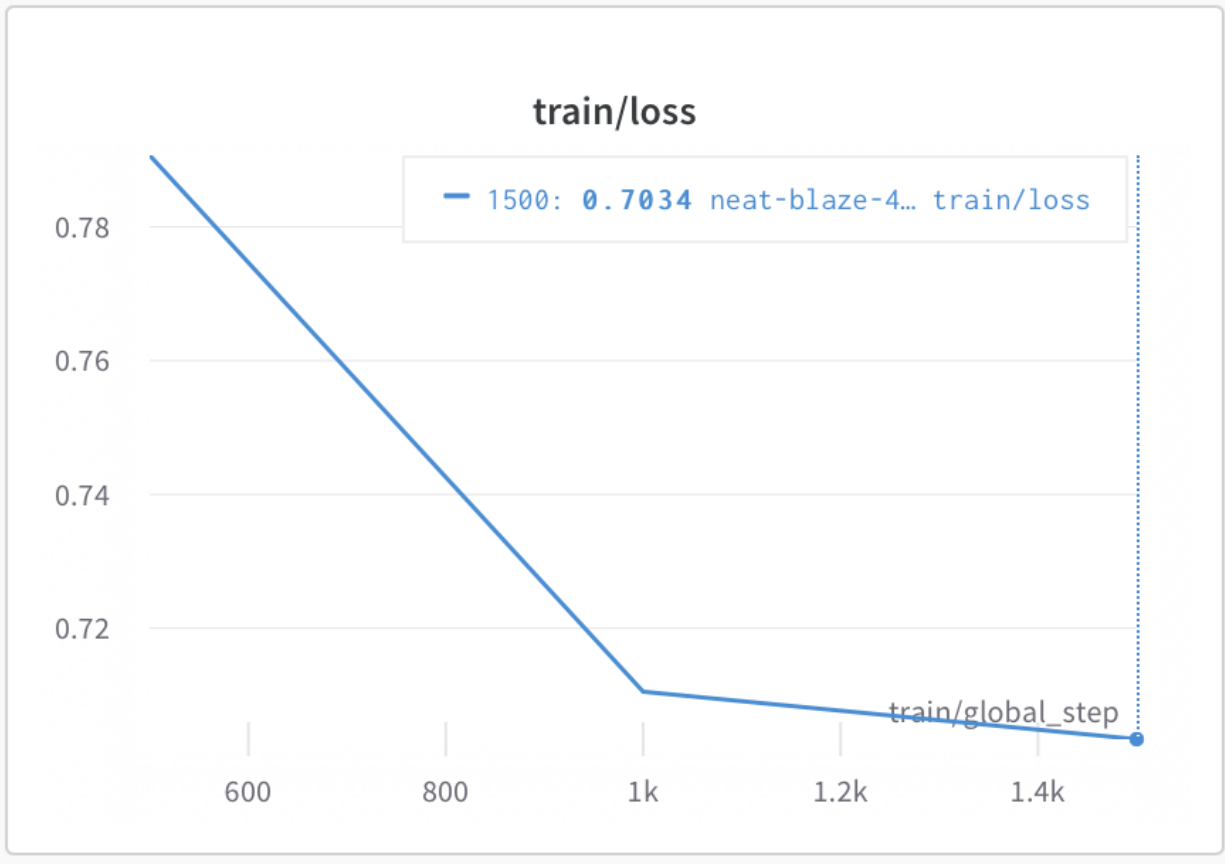
The model is available for download and can be used with various deep learning frameworks such as PyTorch and TensorFlow. It also comes with a user-friendly interface that allows users to input prompts and generate jokes on the fly. Overall, the model is a valuable resource for anyone looking to incorporate humor into their natural language processing applications.
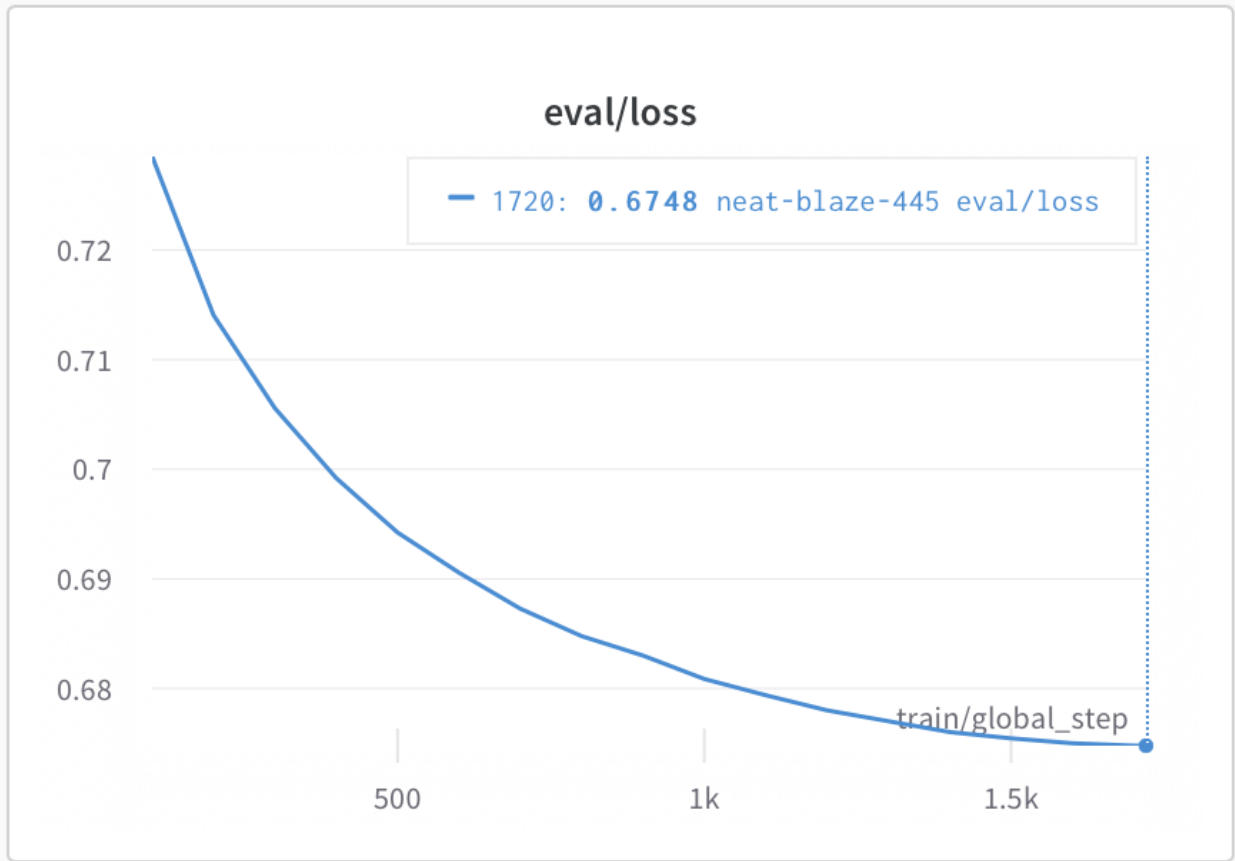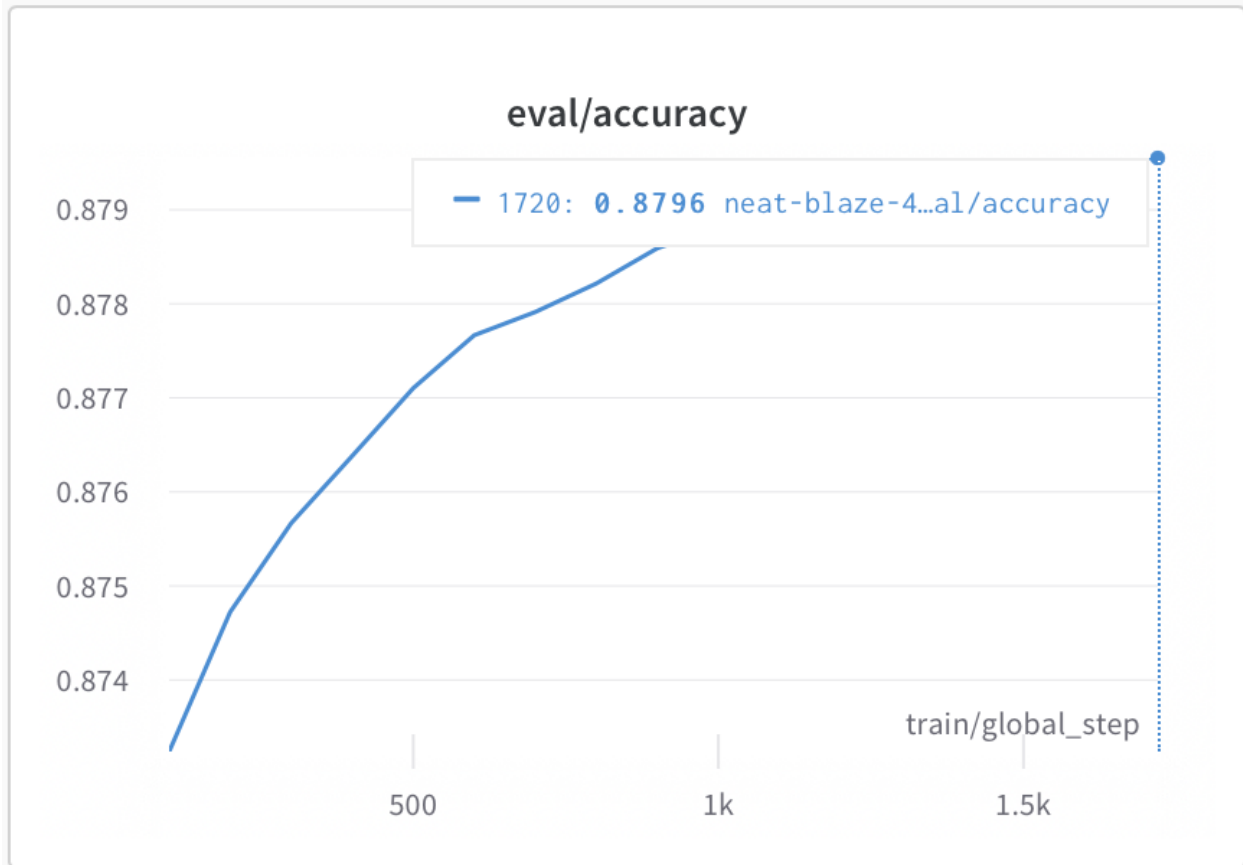
# Training

Ww are use a <u>configuration file for a Kubernetes</u> deployment of a joke generation bot using natural language processing (NLP) techniques. The file specifies the necessary resources and settings for the deployment, such as the container image, the number of replicas, the resource limits, and the networking configuration

For training model we are use the <u>script</u> which loads a pre-trained model and generates jokes based on input prompts or queries. The model is based on the GPT-2 architecture and is fine-tuned on a dataset of jokes to generate humorous content. The script also includes options for specifying the length of the generated jokes and the number of samples to generate.

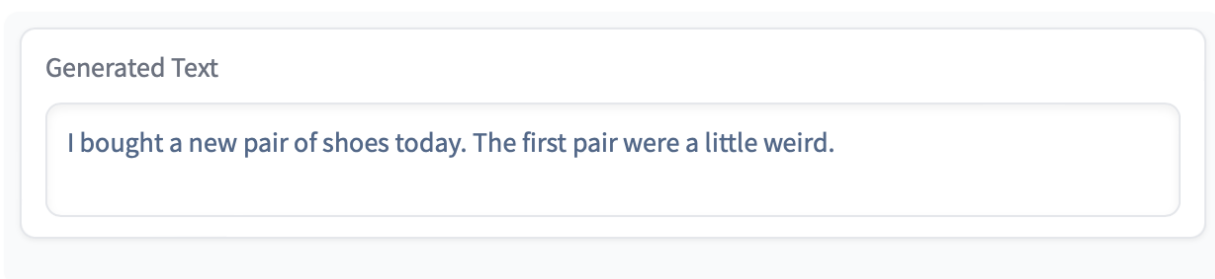Model training metrics can be found at the <u>link</u>. A few screenshots below.

train/loss

1500: **0.7034** neat-blaze-4… train/loss

**eval/loss**

— 1720: **0.6748** neat-blaze-445 eval/loss

## Results

At the moment, our model generates jokes, an example of one of them is below.

Generated Text

> I bought a new pair of shoes today. The first pair were a little weird.

The model has several limitations due to poor data quality and a small model.

## Resources

Fraser, A. (n.d.). Short Jokes Dataset. Hugging Face Datasets. Retrieved March 29, 2023, from **https://huggingface.co/datasets/Fraser/short-jokes**

Korshuk, A. (n.d.). GPT2-jokes. Hugging Face. Retrieved March 29, 2023, from **https://huggingface.co/AlekseyKorshuk/gpt2-jokes**

Korshuk, A. (n.d.). Hugging Face - d8dg1f9i. Weights & Biases. Retrieved March 29, 2023, from **https://wandb.ai/aleksey-korshuk/huggingface/runs/d8dg1f9i?workspace=user-**

# Each Team Member Contribution

- Aliaksei Korshuk - model training

- Rafail Venediktov - search and refinement of the model

- Panov Evgenii - search and analysis of dataset

Each member of the team has made an integral contribution to obtaining this result.

# Plan for the next three weeks

- collect a better dataset with higher quality

- train the bigger model, ex. gpt-neo / gpt-j