

TAG ME UP LAST.FM

Multi-class lyrics classification, a Large Language Model approach

Sangeeths
Chandrakumar

Fachhochschule
Graubünden
Switzerland

sangeeths.chandrakumar@stud.fhgr.ch
mar@stud.fhgr.ch
ORCID 0009-0007-0403-4661

Florian
Klessascheck

Fachhochschule
Graubünden
Switzerland

florian.klessascheck@stud.fhgr.ch
@stud.fhgr.ch
ORCID 0009-0000-0709-7974

Adrian
Joost

Fachhochschule
Graubünden
Switzerland

adrian.joost@stud.fhgr.ch
gr.ch
ORCID 0009-0002-0950-0119

Ana
Petrus

Fachhochschule
Graubünden /
scivia GmbH
Switzerland

ana.petrus@fhgr.ch
ORCID 0000-0002-0928-8894

Abstract – The automatic classification of music tracks according to their lyrics represents an innovative approach for music streaming services. Tagging systems assist users in discovering new music that aligns with their preferences, thereby enhancing satisfaction with the service. Previous research has employed residual neural networks for the classification of music based on their spectrogram. This paper investigates the potential of classifying tracks based solely on their lyrics. To this end, 48'000 tracks with their lyrics and tags were extracted from genius.com and Last.fm. From this dataset, 429 distinct tags were identified for evaluation purposes. Using this dataset, a Mistral-7b-instruct-v2 model was trained and evaluated, demonstrating classification scores of up to 80%. The results indicate that lyrics can serve as a reliable indicator for certain tags.

Keywords – Large Language Models, classification, tagging, music, cataloguing

I. INTRODUCTION

Last.fm is a website that enables users to add tags via a process known as collaborative tagging [1]. Those tags include, but are not limited to, genre, language and the overall feeling (or "vibe") of the song. For the purposes of this research, the lyrics and their respective tags of 48'000 songs from popular (2024) artists were extracted from genius.com and last.fm.

II. STATE OF RESEARCH

The automatic tagging of audio sequences based on machine learning is a well-researched task [2]. For example, in their work "Audio tagging with noisy labels and minimal supervision" [3] Fonseca et al. focus on classifying urban sounds sources and music genres, using spectrum analysis in combination with recurrent neural networks.

Newer research suggests [4] that using lyrics and natural language processing yields better results on genre and "vibe" tagging of music.

III. RESEARCH QUESTION AND METHODOLOGY

The aforementioned works do not employ the use of large language models and are limited to a single prediction for the classification of a single lyric text. This paper builds upon the concept of tagging with the aid of large language models. Consequently, it seeks to answer the following question:

How well can a fine-tuned large language model perform a multi-class audio tag classification based on lyrics and the primary artist's name alone?

The primary artist's name was included to reflect a real-world scenario in which the artist responsible for a particular song is known. From the dataset that we have compiled, we have manually selected 429 of the approximately 4'000 tags that are present. The focus of this research is exclusively on the aforementioned tags. Any other tags will be disregarded, even if they are



present in the predictions. For the purposes of curation, tags that had been observed less than three times across the entire dataset, or that were of low quality (for example, labelled as "good" or "trash"), were filtered out. Any sub-genres or duplicate tags present in a single observation were merged into a single tag.

With regard to the prediction of tags, the mistral-7b-instruct-v2 model was selected, on the grounds of its excellent instruction cohesion and the availability of information on the fine tuning [5] of mistral models.

The initial dataset was divided into three subsets: a training set comprising 66% of the data, a test set comprising 17% of the data, and an evaluation set comprising 17% of the data. The fine-tuning was conducted in accordance with a guide [5] utilising Parameter-Efficient Fine-Tuning (PEFT) [6] with Low-Rank Adaptation (LoRA) [7], with the objective of reducing the computational resource requirements. Furthermore, in order to limit resource usage, all lyrics in the train and test splits exceeding 2'000 tokens were excluded from the data set, resulting in the exclusion of 34 songs (approximately 0.085%) in total.

To assess the impact of fine-tuning, two baselines were established using the unmodified base model. One baseline was provided with all curated tags, while the other was conducted blind, where the model was only presented with the identical prompt as the fine-tuned model:

```
"[INST]<SYS>Tag the song based on the lyrics, only respond in json {"tags": []}</SYS>\n<artist name>\n<lyrics>[/INST]"
```

Due to a lack of resources, two evaluations were conducted. The first evaluation used a fine-tuned model and included 24'000 songs. The second evaluation used the previously described 8'000 songs and compared the fine-tuned model to the base model. However, due to an oversight in data preparation, songs without any tags after filtering were excluded from the evaluation set. This resulted in a reduction of 1'945 songs, leaving 22'055 songs in the first evaluation set.

Confusion scores were calculated on a per-tag basis. For each song, all curated tags were checked and scored in accordance with the following scheme.

A. Evaluation Strategy

This paper asserts the quality of the model by its precision, because the focus of classification is that a predicted tag is correct, not that the model finds all tags.

As the TN (true negative) count is often significantly higher in multi-class classification, the evaluation focused instead on precision, recall and F1 score.

Table 1: Confusion Scoring Strategy

	Tag in label	Tag in prediction
True positive	True	True
False positive	False	True
False negative	True	False
True negative	False	False

IV. RESULTS

Accuracy and specificity were included but as shown in Appendix Figure 1, were heavily influenced by the high TN count. Based on this, looking at the tags ranked by precision, Appendix Figure 2 suggests that the genres country, rap and hip-hop can be well identified by their lyrics. Tags such as electronic, rhythm and blues or trap exhibited lower predictive accuracy.

For the overall best tags (Appendix Figure 3), the model performs well in predicting the languages (Appendix Figure 4) of the songs. The more important precision score is approximately 70% for the top 20 tags in general and by language. Some false positive and false negative songs¹ can be attributed to incorrect or incomplete last.fm tags. In these cases, the model correctly predicted the tag, but it was not included in the labels. For example, there are songs² without German lyrics³ that are wrongly tagged as "German".

A. Model Performance

Figure 1 illustrates that the fine-tuned model exhibits the most optimal performance overall. The elevated precision scores can be attributed to the exclusion of any non-curated tags, even when accounting for false positives. Notably, the baseline that was provided with all available tags demonstrated a less favourable outcome than the baseline that was not assisted.

Run	Name	F1	Recall	Precision
finetuned	Top 20	63.14%	65.28%	61.12%
	Count ≥ 500	64.84%	67.98%	61.97%
	Count ≥ 250	63.83%	67.02%	60.93%
	Total	50.95%	46.14%	56.89%
baseline-tags	Top 20	38.67%	28.83%	58.72%
	Count ≥ 500	41.30%	31.42%	60.23%
	Count ≥ 250	39.61%	29.81%	59.02%
	Total	26.75%	21.73%	34.79%
baseline-blind	Top 20	41.91%	30.89%	65.13%
	Count ≥ 500	46.59%	35.56%	67.52%
	Count ≥ 250	42.60%	31.54%	65.60%
	Total	27.99%	19.81%	47.68%

Figure 1: Comparison of approaches

¹ out of scope of this paper, see https://github.com/Aeolin/tag-me-up-last-fm/blob/master/evaluation_results/confusion/cm_record_including_songs.json for further investigation

² <https://www.last.fm/music/Michael+Jackson/+Happy+Birth+day+Lisa>

³ <https://genius.com/Michael-jackson-happy-birthday-lisa-lyrics>

V. DISCUSSION

Analysis of the results shows that certain tags, such as "rap", are very well suited for evaluation by large language models, while others, such as "electronic", are not. The reason may be that some, but not all, tags have a correlation with the lyrics. Considering this, we suggest that a combination of different predictors or approaches might yield even better results for tagging songs, for example natural language processing for lyrics-heavy tags and sound waves for melody-focused tags.

The baseline model with all available tags performed worse than the blind baseline model, which may stem from an overfilled prompt containing a list of all 429 tags.

During the evaluation, after training, we discovered incomplete and incorrect last.fm tags. It is therefore recommended that, when recreating the approach described in this paper, steps are taken to improve the quality of the dataset by filtering out any unwanted tags or songs with poor label quality.

However, the desired tags must be defined in accordance with the requirements. Given the large number of songs by the same artist in the dataset, it would be beneficial to use a more diverse dataset in terms of artist count, or to use only lyrics, in order to prevent overfitting due to the recognition of the artist's name.

REFERENCES

- [1] Y.-X. Chen, S. Boring, and A. Butz, 'How last. fm illustrates the musical world: user behavior and relevant user-generated content.', in *Proceedings of the international workshop on visual interfaces to the social and semantic web*, Hong Kong, China, 2010. [Online]. Available: <https://www.medien.fki.uni.de/pubdb/publications/pub/chen2010VISSW2/chen2010VISSW2.pdf>
- [2] J. Breebaart and M. F. McKinney, 'Features for Audio Classification', in *Algorithms in Ambient Intelligence*, vol. 2, W. F. J. Verhaegh, E. Aarts, and J. Korst, Eds., in Philips Research, vol. 2, Dordrecht: Springer Netherlands, 2004, pp. 113–129. doi: [10.1007/978-94-017-0703-9_6](https://doi.org/10.1007/978-94-017-0703-9_6).
- [3] E. Fonseca, M. Plakal, F. Font, D. P. W. Ellis, and X. Serra, 'Audio tagging with noisy labels and minimal supervision', 2019, *arXiv*. doi: [10.48550/ARXIV.1906.02975](https://doi.org/10.48550/ARXIV.1906.02975).
- [4] A. Tsaptsinos, 'Lyrics-Based Music Genre Classification Using a Hierarchical Attention Network', 2017, *arXiv*. doi: [10.48550/ARXIV.1707.04678](https://doi.org/10.48550/ARXIV.1707.04678).
- [5] H. Carroll, I. Dhanani, N. Khalil, V. Parashar, T. Fong, and Anarkoic, 'notebooks/mistral-finetune-own-data.ipynb at main · brevdev/notebooks', Features for Audio and Music Classification. Accessed: Apr.

24, 2024. [Online]. Available:

<https://github.com/brevdev/notebooks/blob/main/mistral-finetune-own-data.ipynb>

- [6] L. Xu, H. Xie, S.-Z. J. Qin, X. Tao, and F. L. Wang, 'Parameter-Efficient Fine-Tuning Methods for Pretrained Language Models: A Critical Review and Assessment', 2023, *arXiv*. doi: [10.48550/ARXIV.2312.12148](https://doi.org/10.48550/ARXIV.2312.12148).
- [7] Y. Yu *et al.*, 'Low-rank Adaptation of Large Language Model Rescoring for Parameter-Efficient Speech Recognition', in *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2023, pp. 1–8. doi: [10.1109/ASRU57964.2023.10389632](https://doi.org/10.1109/ASRU57964.2023.10389632).

ACKNOWLEDGEMENTS

This study was conducted as part of the course "CDS1091 Natural Language Processing", taught by Prof. Corsin Capol, in the context of BSc Computational and Data Science at the University of Applied Sciences of the Grisons (Fachhochschule Graubünden), Switzerland.

CONTRIBUTIONS

Sangeeths Chandrakumar: conceptualisation, data curation, formal analysis, investigation, methodology, software, validation, visualisation, writing – original draft

Florian Klessascheck: conceptualisation, data curation, formal analysis, investigation, methodology, software, validation, visualization, writing – original draft

Adrian Joost: conceptualisation, writing – original draft

Ana Petrus: supervision, writing – review & editing

APPENDIX

Appendix Figures 1 through 4.



		Tag / Predicted Count / Actual Count																		
	country	hard rock	rap	hip-hop	jazz	christian	soul	alternative rock	dance	pop	80s	metal	rock	electronic	indie pop	rhythm and blues	trap	folk	indie	alternative
Precision	652 692	281 648	7297 5626	7415 5770	367 662	232 627	312 1507	893 1127	326 1627	6777 5326	803 954	1229 971	4082 2974	1496 1489	1290 1116	1798 1397	2595 1783	1494 817	1461 1908	2562 1938
	87.9%	83.9%	74.3%	73.9%	72.9%	69.3%	67.1%	65.0%	64.8%	63.7%	61.3%	57.7%	54.1%	51.8%	51.5%	51.3%	47.0%	46.8%	42.0%	33.4%
Recall	79.5%	35.5%	92.4%	91.0%	33.4%	25.2%	12.7%	51.0%	12.7%	76.9%	49.6%	71.3%	71.0%	50.2%	55.6%	62.8%	65.2%	78.6%	30.2%	41.4%
	71.6%	33.2%	70.0%	68.9%	29.7%	22.7%	12.0%	40.0%	11.8%	53.5%	37.7%	46.8%	44.3%	34.2%	36.5%	39.4%	37.6%	41.5%	21.3%	22.7%
F1																				

Appendix Figure 2: Scores for top 20 tags by actual count, ordered by precision

	Tag / Predicted Count / Actual Count																			
	hip-hop	rap	pop	rock	alternative	indie	trap	dance	soul	electronic	rhythm and blues	alternative rock	indie pop	metal	80s	folk	country	jazz	hard rock	christian
Accuracy	7415 5770	7297 5626	6777 5326	4082 2974	2562 1938	1461 1908	2595 1783	326 1627	312 1507	1486 1489	1798 1397	893 1127	1290 1116	1229 971	803 954	1494 817	652 692	367 662	281 648	232 627
	84.6%	85.3%	79.5%	83.4%	83.1%	85.7%	86.5%	88.2%	88.8%	88.7%	89.0%	91.1%	90.2%	91.5%	91.5%	91.0%	93.9%	92.6%	92.8%	92.5%
Specificity	88.6%	89.1%	86.1%	90.6%	92.1%	96.1%	93.5%	99.5%	99.5%	96.6%	96.0%	98.5%	97.2%	97.6%	98.6%	96.6%	99.6%	99.6%	99.8%	99.7%
	73.9%	74.3%	63.7%	54.1%	33.4%	42.0%	47.0%	64.8%	67.1%	51.8%	51.3%	65.0%	51.5%	57.7%	61.3%	46.9%	87.9%	72.9%	83.9%	69.3%
Precision	73.9%	74.3%	63.7%	54.1%	33.4%	42.0%	47.0%	64.8%	67.1%	51.8%	51.3%	65.0%	51.5%	57.7%	61.3%	46.9%	87.9%	72.9%	83.9%	69.3%
	91.0%	92.4%	76.9%	71.0%	41.4%	30.2%	65.2%	12.7%	12.7%	50.2%	62.8%	51.0%	55.6%	71.3%	49.6%	78.6%	79.5%	33.4%	35.5%	25.2%
Recall	91.0%	92.4%	76.9%	71.0%	41.4%	30.2%	65.2%	12.7%	12.7%	50.2%	62.8%	51.0%	55.6%	71.3%	49.6%	78.6%	79.5%	33.4%	35.5%	25.2%
	68.9%	70.1%	53.5%	44.3%	22.7%	21.3%	37.6%	11.8%	12.0%	34.2%	39.4%	40.0%	36.5%	46.8%	37.7%	41.6%	71.6%	29.7%	33.2%	22.7%
F1	68.9%	70.1%	53.5%	44.3%	22.7%	21.3%	37.6%	11.8%	12.0%	34.2%	39.4%	40.0%	36.5%	46.8%	37.7%	41.6%	71.6%	29.7%	33.2%	22.7%

Appendix Figure 1: Scores for top 20 tags by actual count, ordered by the actual count



	Tag / Predicted Count / Actual Count																			
	romanian	arabic	k-pop	norwegian	turkish	metal	korean	indonesian	j-pop	spanish	swedish	azerbaijan	rusia	italian	french rap	serbian	german rap	british	danish	japanese
	190	47	134	61	66	1229	77	118	174	294	258	43	158	171	354	31	62	282	27	165
	194	46	130	63	56	971	83	60	125	246	434	122	252	189	511	90	143	570	112	211
F1	93.4%	87.8%	71.1%	64.0%	53.2%	46.8%	44.5%	43.6%	43.3%	37.3%	37.0%	35.2%	33.6%	33.1%	29.9%	28.7%	26.7%	25.1%	24.1%	22.5%
Precision	97.9%	93.5%	83.1%	80.0%	66.1%	57.7%	64.5%	47.2%	52.0%	51.3%	75.0%	100.0%	65.2%	54.0%	57.1%	87.1%	70.5%	63.2%	100.0%	41.8%
	95.4%	93.5%	83.1%	76.2%	73.2%	71.3%	59.0%	85.0%	72.0%	57.7%	42.2%	35.2%	40.9%	46.0%	38.6%	30.0%	30.1%	29.5%	24.1%	32.7%

Appendix Figure 4: Scores for top 20 language related tags by F1

	Tag / Predicted Count / Actual Count																			
	romanian	grunge	grindcore	industrial	country	k-pop	rap	hip-hop	heavy metal	nu metal	christian rock	60s	pop	hippy	thrash metal	metal	rock	j-pop	gospel	punk rock
	190	433	245	304	652	134	7297	7415	483	278	212	148	6777	125	372	1229	4082	174	218	315
	194	409	238	309	692	130	5626	5770	570	289	152	152	5326	167	539	971	2974	125	450	382
F1	93.4%	90.7%	88.4%	74.3%	71.6%	71.1%	70.0%	68.9%	63.6%	63.5%	62.5%	57.4%	53.5%	53.3%	51.2%	46.8%	44.3%	43.3%	43.2%	41.9%
Precision	97.9%	93.0%	94.5%	88.2%	87.9%	83.1%	74.3%	73.9%	85.8%	79.6%	67.8%	77.2%	63.7%	85.2%	84.7%	57.7%	54.1%	52.0%	93.1%	66.6%
	95.4%	97.3%	93.3%	82.5%	79.5%	83.1%	92.4%	91.0%	71.1%	75.8%	88.8%	69.1%	76.9%	58.7%	56.4%	71.3%	71.0%	72.0%	44.7%	53.1%

Appendix Figure 3: Scores for top 20 tags by F1 and > 100 actual counts, ordered by F1

