# MAEST: Open music representation n for music understanding



<u>Pablo Alonso-Jiménez, Xavier Serra, Dmitry Boqdanov</u>

Music Technology Group, Universitat Pompeu Fabra

# Motivation

Descriptive tags are **difficult to obtain** and **noisy**. We need alternative ways of generating training targets for large music collections and suitable training approaches to

# **Experiments and results**

Extracting embeddings from the transformer

#### develop music representation models.

#### U - 36.5 38.8 40.2 40.7 40.4 40.1 39.3 38.8 -40 ated -38.1 40.1 40.6 40.4 40.2 39.3 39.6 39.4 b - 39 aten 39.6 40.9 41.0 40.7 40.2 39.3 39.3 - 38 ЫС - 39.2 40.5 41.1 40.6 40.7 39.5 39.4 39.5 ů cda - 37 40.6 41.3 41.1 40.6 39.9 39.4 39.4 39.1 12 Transformer block

# **Experimental setup**



### Impact of the initial weights

| Model                    | RW   | DeiT | PaSST |
|--------------------------|------|------|-------|
| Pre-training task: Disco | gs20 |      |       |
| MAEST-10s                | 20.5 | 22.7 | 22.8  |
| MAEST-10s-swa            | 20.1 | 23.2 | 23.5  |
| Downstream task: MTT     |      |      |       |
| MAEST-10s                | 38.7 | 40.4 | 41.1  |
| MAEST-10s-swa            | 39.0 | 40.2 | 41.0  |



# Conclusions

1. Patchout allows for **efficient training** and **inference** with Transformers

### Effect of the input segment length

| Model                  | 5s     | 10s  | 20s  | 30s  |
|------------------------|--------|------|------|------|
| Pre-training task: Dis | cogs20 |      |      |      |
| MAEST-T                | 21.1   | 22.8 | 24.8 | 26.1 |
| MAEST-T-swa            | 21.3   | 23.5 | 25.8 | 27.0 |
| Downstream task: MT    | Т      |      |      |      |
| MAEST-T                | 40.8   | 41.1 | 41.2 | 41.7 |
| MAEST-T-swa            | 40.9   | 41.0 | 41.2 | 41.5 |

#### Faster feature extraction with inference patchout



2. Transformers allow for **better music representations** than CNNs

3. We propose MAEST, an publicly available music representation model

throughput (analyzed minutes / second)

#### **Contact**: <u>pablo.alonso@upf.edu</u> @pablo alonso

Park, J. Lee, J.W. Ha, and J. Nam. "Representation learning of music using artist label.," in Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR). 2018. [1]

Saeed, Aaqib, Grangier, David, and Zeghidour, Neil. "Contrastive learning of general-purpose audio representations." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2021. [2]

Tan, Mingxing, and Quoc Le. "Efficientnet: Rethinking model scaling for convolutional neural networks." International Conference on Machine Learning (ICML). 2019. [3]

Hershey, Shawn, et al. "CNN architectures for large-scale audio classification." International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2017. [4]



# MAEST: Open music representation n for music understanding



Pablo Alonso-Jiménez, Xavier Serra, Dmitry Bogdanov

Music Technology Group, Universitat Pompeu Fabra

# Motivation

Descriptive tags are **difficult to obtain** and **noisy**. We need alternative ways of generating training targets for large music collections and suitable training approaches to

# **Experiments and results**

Extracting embeddings from the transformer

#### develop **music representation** models.

#### -40 ated - 38.1 40.1 40.6 40.4 40.2 39.3 39.6 39.4 b - 39 aten 39.6 40.9 41.0 40.7 40.2 39.3 39.3 - 38 ЫС - 39.2 40.5 41.1 40.6 40.7 39.5 39.4 39.5 0 cda - 37 39.1 40.6 41.3 41.1 40.6 39.9 39.4 39.4 Transformer block

# **Experimental setup**



## Impact of the initial weights Effect of the input segment length

|                          |      |      |       | 2                      |        |      |      |      |  |
|--------------------------|------|------|-------|------------------------|--------|------|------|------|--|
| Model                    | RW   | DeiT | PaSST | Model                  | 5s     | 10s  | 20s  | 30s  |  |
| Pre-training task: Disco | gs20 |      |       | Pre-training task: Dis | cogs20 |      |      |      |  |
| MAEST-10s                | 20.5 | 22.7 | 22.8  | MAEST-T                | 21.1   | 22.8 | 24.8 | 26.1 |  |
| MAEST-10s-swa            | 20.1 | 23.2 | 23.5  | MAEST-T-swa            | 21.3   | 23.5 | 25.8 | 27.0 |  |
| Downstream task: MTT     |      |      |       | Downstream task: MT    | Т      |      |      |      |  |
| MAEST-10s                | 38.7 | 40.4 | 41.1  | MAEST-T                | 40.8   | 41.1 | 41.2 | 41.7 |  |
| MAEST-10s-swa            | 39.0 | 40.2 | 41.0  | MAEST-T-swa            | 40.9   | 41.0 | 41.2 | 41.5 |  |
|                          |      |      |       |                        |        |      |      |      |  |



# Conclusions

1. Patchout allows for **efficient training** and **inference** with Transformers

#### Faster feature extraction with inference patchout



2. Transformers allow for **better music representations** than CNNs

3. We propose MAEST, an publicly available music representation model

throughput (analyzed minutes / second)

#### 

[1] Park, J. Lee, J.W. Ha, and J. Nam. "Representation learning of music using artist label.," in Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR). 2018.

[2] Saeed, Aaqib, Grangier, David, and Zeghidour, Neil. "Contrastive learning of general-purpose audio representations." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2021.

[3] Tan, Mingxing, and Quoc Le. "Efficientnet: Rethinking model scaling for convolutional neural networks." International Conference on Machine Learning (ICML). 2019.

[4] Hershey, Shawn, et al. "CNN architectures for large-scale audio classification." International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2017.



# MAEST: Music Audio Efficient Spectrogram Transformer



## **Efficient Supervised Training of Audio Transformers for Music Representation Learning**

Pablo Alonso-Jiménez, Xavier Serra, Dmitry Bogdanov

Music Technology Group, Universitat Pompeu Fabra

# Motivation

Our goal is to propose **music representation models** with a focus in **semantic music description**. We rely on **convolution-free Transformers**, and propose experiments

## Impact of the initial weights

| Model                    | RW    | DeiT | PaSST |
|--------------------------|-------|------|-------|
| Pre-training task: Disco | ogs20 |      |       |
| MAEST-10s                | 20.5  | 22.7 | 22.8  |
| MAEST-10s-swa            | 20.1  | 23.2 | 23.5  |

to understand which conditions **optimize** the **downstream performance**.

| Exper | 'imenta | setup |
|-------|---------|-------|
|       |         |       |

Training



Downstream evaluation



| Downstream task: MTT |      |      |      |
|----------------------|------|------|------|
| MAEST-10s            | 38.7 | 40.4 | 41.1 |
| MAEST-10s-swa        | 39.0 | 40.2 | 41.0 |

Starting the training from the PaSST pre-trained weights produces the best performance in the downstream task.

#### Effect of the input segment length

| Model                  | 5s     | 10s  | 20s  | 30s  |
|------------------------|--------|------|------|------|
| Pre-training task: Dis | cogs20 |      |      |      |
| MAEST-T                | 21.1   | 22.8 | 24.8 | 26.1 |
| MAEST-T-swa            | 21.3   | 23.5 | 25.8 | 27.0 |
| Downstream task: MT    | T      |      |      |      |
| MAEST-T                | 40.8   | 41.1 | 41.2 | 41.7 |
| MAEST-T-swa            | 40.9   | 41.0 | 41.2 | 41.5 |

We modify the input sequence length, from 5 up to 30 seconds finding that the results **consistently increase** in the downstream task.

### Performance in downstream tasks

|   | MTGJ                        | -Genre                      | MTG                         | J-Inst                      | MTGJ-                       | Mood                        | MTG.                        | J-T50                              | MT                          | AT                          | MS                          | SDs                         | MS                          | Dc                          |
|---|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|------------------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
|   | ROC                         | mAP                         | ROC                         | mAP                         | ROC                         | mAP                         | ROC                         | mAP                                | ROC                         | mAP                         | ROC                         | mAP                         | ROC                         | mAP                         |
| State of the art                                  |                             |                             |                             |                             |                             |                             |                             |                                    |                             |                             |                             |                             |                             |                             |
| Fully-trained                                     | -                           | -                           | -                           | -                           | 77.8<br>[42]                | 15.6<br>[42]                | 83.2<br>[34]                | 29.8<br>[34]                       | 90.69<br>[41]               | 38.44<br><i>[41]</i>        | 92.2<br>[40]                | 38.9<br>[40]                | 89.7<br>[40]                | 34.8<br>[40]                |
| Embeddings  | 87.7<br>[6]                 | 19.9<br>[6]                 | 77.6<br>[6]                 | 19.8<br>[6]                 | 78.6<br>[5] <sup>†</sup>    | 16.1<br>[5] <sup>†</sup>    | 84.3<br>[5] <sup>†</sup>    | 32.1<br>[5] <sup>†</sup>           | 92.7<br>[7] <sup>†</sup>    | 41.4<br>[5] <sup>†</sup>    | -                           | -                           | 90.3<br>[5] <sup>†</sup>    | 36.3<br>[5] †               |
| Baseline<br>EffNet-B0                             | 87.7                        | 19.9                        | 77.6                        | 19.8                        | 75.6                        | 13.6                        | 83.1                        | 29.7                               | 90.2                        | 37.4                        | 90.4                        | 32.8                        | 88.9                        | 32.8                        |
| Our models<br>MAEST-10s<br>MAEST-20s<br>MAEST-30s | 88.1<br>88.1<br><b>88.2</b> | 21.1<br>21.4<br><b>21.6</b> | 79.7<br>79.9<br><b>80.0</b> | 22.4<br>22.6<br><b>22.9</b> | 77.9<br>77.9<br><b>78.1</b> | 15.1<br>15.2<br><b>15.4</b> | 84.0<br><b>84.1</b><br>84.0 | 31.3<br><b>31.5</b><br><b>31.5</b> | 91.8<br>91.8<br><b>92.0</b> | 41.0<br>41.0<br><b>41.9</b> | 91.5<br>92.1<br><b>92.4</b> | 36.9<br>39.2<br><b>40.7</b> | 88.9<br>89.5<br><b>89.8</b> | 32.7<br>34.5<br><b>35.4</b> |



 we evaluate our models in the MTG Jamendo Dataset, MagnaTagAtune and the Million Song Dataset.

\* We find that our models are the **best performing open solution**.

#### Faster feature extraction with inference patchout



We experiment applying patchout at inference time.

It is possible to reach scenarios where the throughput is higher than

# **Experiments and results**

## Extracting embeddings from the transformer



We experimented with the class c), distillation d), and average a) tokens. Staking the three tokens produces the best performance.
The optimal features are in the middle blocks of the transformer.

the fully convolutional baselines while keeping higher performance.

# Conclusions

- Patchout allows for efficient training and inference with Transformers.
- Transformers allow for better music representations than CNNs
- We propose MAEST, an publicly available music representation model.

## Contact: <u>pablo.alonso@upf.edu</u> <u>@pablo\_alonso</u>



# **MAEST: Open music representation n** for music understanding



Pablo Alonso-Jiménez, Xavier Serra, Dmitry Bogdanov

Music Technology Group, Universitat Pompeu Fabra

# Motivation

Descriptive tags are difficult to obtain and noisy. We need alternative ways of generating training targets for large music collections and suitable training approaches to develop **music representation** models.

## **Experimental setup**





## **Experiments and results**

**Extracting embeddings from the transformer** 



#### Impact of the initial weights

| Model                    | RW    | DeiT | PaSST |
|--------------------------|-------|------|-------|
| Pre-training task: Disco | ogs20 |      |       |
| MAEST-10s                | 20.5  | 22.7 | 22.8  |
| MAEST-10s-swa            | 20.1  | 23.2 | 23.5  |
| Downstream task: MTT     |       |      |       |
| MAEST-10s                | 38.7  | 40.4 | 41.1  |
| MAEST-10s-swa            | 39.0  | 40.2 | 41.0  |

#### **Effect of the input segment length**

| Model                 | 5s      | 10s  | 20s  | 30s  |
|-----------------------|---------|------|------|------|
| Pre-training task: Di | scogs20 |      |      |      |
| MAEST-T               | 21.1    | 22.8 | 24.8 | 26.1 |
| MAEST-T-swa           | 21.3    | 23 5 | 25.8 | 27.0 |

#### Faster feature extraction with inference patchout



# Conclusions

Patchout allows for efficient training and inference with Transformers

2. Transformers allow for **better music representations** than CNNs

3. We propose MAEST, an publicly available music representation model

#### **Contact**: <u>pablo.alonso@upf.edu</u> @pablo alonso

Park, J. Lee, J.W. Ha, and J. Nam. "Representation learning of music using artist label.," in Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR). 2018. [1]

Saeed, Aaqib, Grangier, David, and Zeghidour, Neil. "Contrastive learning of general-purpose audio representations." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2021. [2]

Tan, Mingxing, and Quoc Le. "Efficientnet: Rethinking model scaling for convolutional neural networks." International Conference on Machine Learning (ICML). 2019. [3]

Hershey, Shawn, et al. "CNN architectures for large-scale audio classification." International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2017. [4]

