# LP-MusicCaps: LLM-Based Pseudo Music Captioning

Music and Audio Computing Lab

Neutune

SeungHeon Doh, Keunwoo Choi, Jongpil Lee, Juhan Nam

Graduate School of Culture Technology, KAIST, South Korea Gaudio Lab, Inc., South Korea<sup>2</sup> Neutune, South Korea<sup>3</sup>

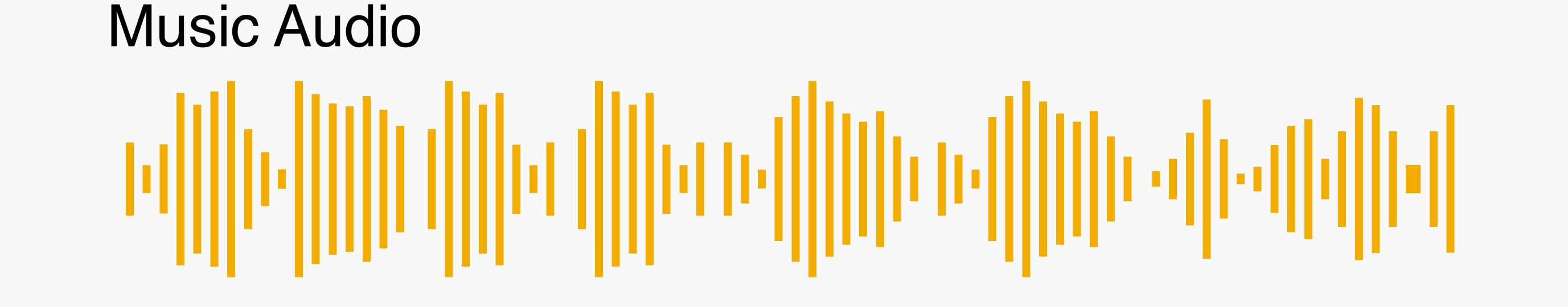
#### Generate Pseudo Caption Using Instruction & Large Language Model Instruction Tag List Write a song description guitar, piano, sentence including the percussion, relaxing melody, slow tempo. following attributes

#### Large Language Model

Generated Pseudo Caption

This gentle guitar song featuring a piano, percussion, and a soothing melody is perfect for relaxation with a slow tempo.

#### Pre-train Music Captioning Model Using Pseudo-Caption Dataset



#### Music Captioning Model (Pretraining)

Generated Pseudo Caption

Methods

Baseline

Writing

Summary

Tag Concat [2, 13]

Proposed Instruction

Template [14]

K2C Aug. [22]

This gentle guitar song featuring a piano, percussion, and a soothing melody is perfect for relaxation with a slow tempo.

**Params** 

220M

175B+

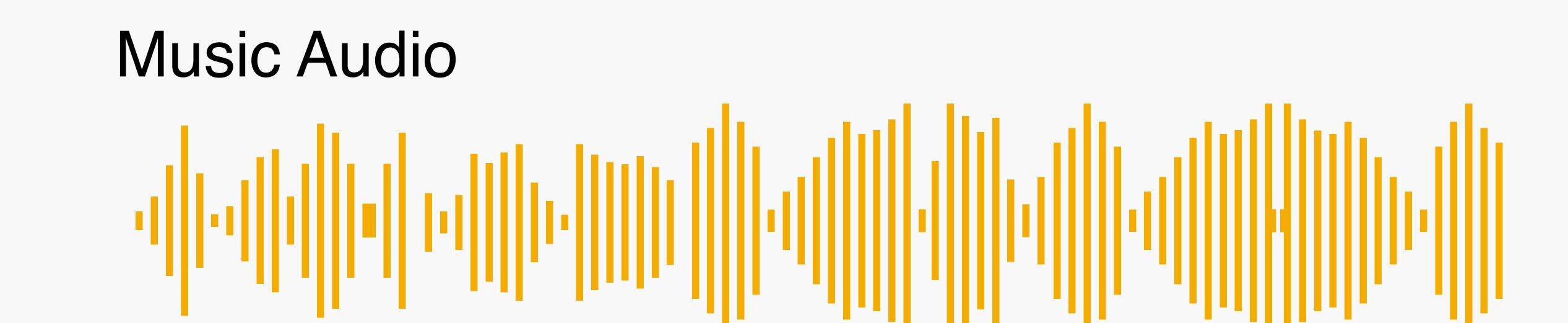
175B+

GPT3.5

GPT3.5

**B**1↑

#### Transfer Learning Using Small **Human-Annotated Caption Dataset**



### Music Captioning Model

(Transfer Learning)

Human Annotated Caption

BERT-S↑

86.24

87.92

86.33

89.26

89.88

R-L

19.52

21.36

**25.83** 

This song features a pedal steel guitar playing the main melody. This starts off with a slide to a high register followed by plucking low register strings. This is accompanied by percussion playing a simple beat in common time....

Diversity Metrics

Vocab<sup>↑</sup>

3506

3507

3760

5521

4198

 $Novel_v \uparrow$ 

46.92

46.93

67.66

56.17

49.52

Length

Avg.Token

 $20.6 \pm 11.2$ 

 $25.6 \pm 11.2$ 

 $14.7 \pm 5.1$ 

 $44.4 \pm 17.3$ 

 $28.6 \pm 10.7$ 

## Part 1. Tag-to-Caption

LLM-based caption generation shows higher performance

LLM can generate various captions depending on the type of instruction and they have different advantage as follows:

• Writing: N-Gram overlapping

• Summary: Semantic similarity

Paraphrase: Diversity

Attribute prediction: Diversity

than other methods.

 $47.9 \pm 18.7$ Paraphrase 88.71 Attribute Prediction 6995  $66.2 \pm 21.6$ **34.09** 88.56 63.16 175B +Table 2. Performance of existing pseudo caption generation methods and the proposed method. LM stand for the language model. Avg. Token stand for the average number of token per caption.

Supervised Metrics

**B**4↑

5.42

6.15

5.52

25.57

27.58

**B**3↑

10.00

1.58

11.37

8.80

16.15

3.01

19.85

14.58



LLM pseudo captions are