



# تولید موسیقی مبتنی بر پردازش زبان طبیعی

محمد امین کیانی ۴۰۰۳۶۱۳۰۵۲





# فهرست مطالب

← ۳- انگیزه و اهمیت موضوع (مقدمه و معرفی)

← ۴- مروری بر روش‌های اولیه

← ۵- جهش به مدل‌های عمیق

← ۶- معیارهای ارزیابی مدل‌ها

← ۷- بررسی معماری و پیاده‌سازی عملی مدل MusicGen

← ۸- چالش‌ها و محدودیت‌ها

← ۹- جمع‌بندی و پیشنهادات (آینده تحقیق)

← ۱۰- مراجع



## • هدف تحقیق:

تولید موسیقی از متن، مسئله‌ای چندرشته‌ای و نوین در مرزهای بین پردازش زبان طبیعی و یادگیری عمیق است که امکان تعامل خلاقانه انسان و ماشین را فراهم می‌سازد.

## • اهمیت موضوع:

- رشد چشمگیر کاربردهای هوش مصنوعی در هنرهای خلاقانه
- نیاز به سیستم‌هایی برای تولید موسیقی بر اساس احساسات یا سناریوهای توصیفی
- تسهیل ساخت موسیقی برای فیلم، بازی، متاورس و تولیدات بدون نیاز به آهنگ‌سازی دستی

## • دلایل انتخاب موضوع:

- نبود راهکارهای دقیق برای ارزیابی موسیقی‌های تولیدشده توسط مدل‌های متنی
- چالش‌های ترکیب درک زبانی با ساختار زمانی و صوتی موسیقی
- رشد سریع مدل‌های ژنراتیو مانند MusicGen، MusicLM، JEN-1 و نیاز به مقایسه مؤثر آن‌ها





## ◆ رویکردهای ابتدایی:

سیستم‌های اولیه مانند (2014) TransProse با تکیه بر قواعد معنایی ساده و واژگان احساسی، نت‌های موسیقی را مستقیماً از تحلیل احساسی متن استخراج می‌کردند.

## ◆ مدل‌های کلاسیک یادگیری ماشین:

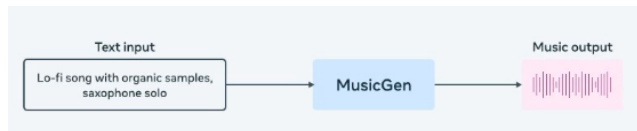
استفاده از مدل‌های RNN و LSTM برای پردازش توالی‌های متنی و تولید نت، با کنترل ویژگی‌هایی مثل شدت و زمان‌بندی صداها.

## ◆ محدودیت‌ها:

- وابستگی به قواعد دست‌نویس
- کیفیت پایین موسیقی تولیدی
- عدم تنوع در خروجی
- عدم توانایی در فهم عمیق مفاهیم متن



◆ با پیشرفت یادگیری عمیق، نسل جدیدی از مدل‌ها با معماری‌های پیچیده ظاهر شد که می‌توانستند روابط پنهان بین زبان و موسیقی را به صورت داده‌محور استخراج کنند.



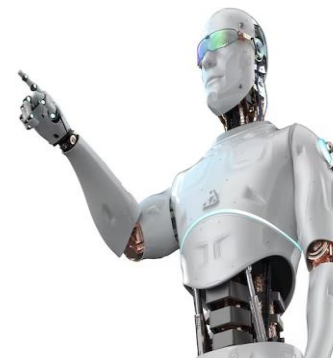
## ◆ مدل‌های شاخص:

- **Riffusion (2022)** : استفاده از مدل دیفیوژن برای تولید طیف‌نگاره‌ها (Spectrogram) از متن.
- **MusicLM (Google, 2023)** : مبتنی بر مدل دو مرحله‌ای رمزگذار-رمزگشا با کدک صوتی و ترنسفورمر.
- **MusicGen (Meta, 2023)** : تولید مستقیم waveform با مدل ترنسفورمری آموزش‌دیده روی داده‌های موسیقی-متن.
- **JEN-1 (2024)** : مدل چندوجهی با توانایی تولید صدای خواننده با کلمات مشخص‌شده.

◆ این مدل‌ها با استفاده از داده‌های عظیم، از ویژگی‌های معنایی و ساختاری متن، موسیقی‌هایی با کیفیت واقعی تولید می‌کنند.

- ارزیابی کیفیت موسیقی تولیدشده، فقط با گوش دادن کافی نیست بلکه به معیارهای کمی و قابل مقایسه نیاز داریم.

- استفاده ترکیبی از این معیارها، تصویری جامع از وفاداری معنایی، کیفیت شنیداری و شباهت آماری می‌دهد.



## معیارهای اصلی:

1. **MOS (Mean Opinion Score)** : میانگین نمره شنوندگان (۵ نمره)، معیاری سنتی ولی ذهنی.

2. **FAD (Fréchet Audio Distance)** : فاصله توزیع ویژگی‌های صدای تولیدی با صدای واقعی که مشابه FID در تصویر است.

3. **CLAP Score (Contrastive Language-Audio Pretraining)** : شباهت معنایی بین متن ورودی و موسیقی خروجی با مدل CLAP.

4. **KAD (Kernel Audio Distance)** : معیار جدید و سبک‌وزن برای اندازه‌گیری اختلاف برداری صوت‌ها.



دانشکده مهندسی کامپیوتر

ارائه‌ی میانی پردازش زبان و گفتار

# MusicGen کلی

- نوع مدل: ترنسفورمر
- ورودی: توصیف متنی آهنگ مورد نظر
- خروجی: دنباله‌ای از بردارهای صوتی که پس از پردازش تبدیل به موج صوتی WAV می‌شوند.

## اجزای اصلی:

### 1. متن پرداز (Text Encoder) :

- برای تبدیل متن به Embedding معنایی.
- این بردار معنایی راهنمای تولید موسیقی است.

### 2. کدک صوتی (Audio Codec) :

- استفاده از EnCodec (مدل کدکننده صوتی از Meta) برای فشرده‌سازی و رمزگشایی صدا.

- صوت نهایی از دنباله‌ی کدهای توصیف‌شده بازسازی می‌شود.

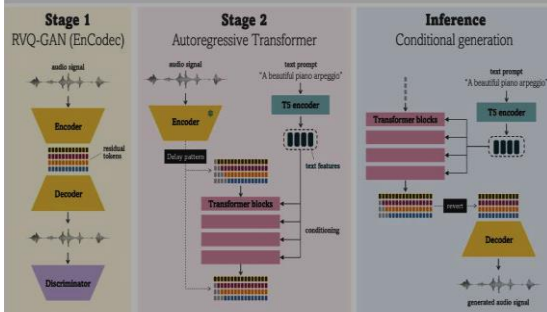
### 3. ترنسفورمر مولد (Music Transformer Decoder) :

- دنباله‌ی کدهای صوتی را از متن رمزگذاری‌شده تولید می‌کند.
- از مکانیزم attention برای هماهنگی بین متن و ساختار زمانی صدا بهره می‌برد.

کیفیت	اندازه پارامتر	نسخه
سریع و سبک	~300M	small
بالانس دقت و سرعت	~1.5B	medium
کیفیت بالا	~3.3B	large
قابلیت کنترل دقیق‌تر	+ ورودی ملودی	melody



## MusicGen



```

# نصب کتابخانه‌های لازم
!pip install transformers datasets scipy

from transformers import AutoProcessor, MusicGenForConditionalGeneration

from IPython.display import Audio

# بارگذاری مدل و پردازشگر
processor = AutoProcessor.from_pretrained("facebook/musicgen-small")

model = MusicGenForConditionalGeneration.from_pretrained("facebook/musicgen-small")

# تعریف پرامپت متنی (می‌تواند فارسی یا انگلیسی باشد)
text_prompt = "A calm Persian classical music with santur and a slow rhythm"

# آماده‌سازی ورودی مدل
inputs = processor(text=[text_prompt], return_tensors="pt")

# تولید توکن‌های صوتی به صورت خونبارگشتی
audio_tokens = model.generate(**inputs, do_sample=True, guidance_scale=3.0,
                               max_new_tokens=256)

# به موج صوتی و پخش آن (tokens) تبدیل خروجی مدل
audio_array = audio_tokens[0].numpy()

# نرخ نمونه‌برداری (مثلاً 32000 هرتز)
sampling_rate = model.config.audio_encoder.sampling_rate

Audio(audio_array, rate=sampling_rate)
  
```

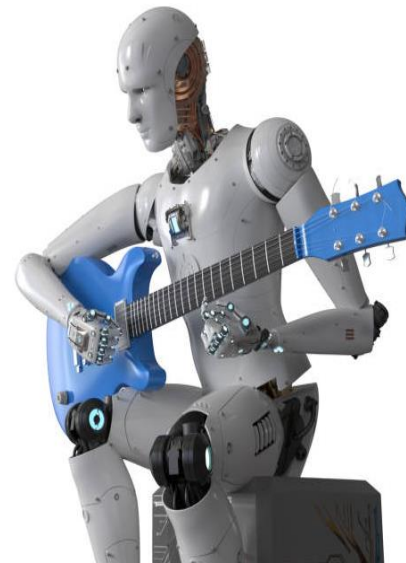


## چالش‌های مفهومی:

- ابهام در زبان طبیعی:  
کلمات می‌توانند بار احساسی یا سبک موسیقایی مبهم یا چندپهلوی داشته باشند.
- نگاشت پیچیده متن به ویژگی‌های صوتی:  
احساسات یا سبک‌ها در متن الزاماً معادل دقیق و مستقیمی در ویژگی‌های آکوستیکی ندارند.
- نبود معیار ارزیابی انسانی:  
بیشتر معیارها آماری و غیرشنیداری‌اند. ارزیابی واقعی کیفیت موسیقی تولیدی نیازمند گوش انسان است.

## چالش‌های پیاده‌سازی:

- مصرف منابع بالا:  
مدل‌های حجیمی مانند MusicLM یا CLAP نیاز به VRAM و RAM بالا دارند. اجرا روی سیستم‌های ضعیف یا Google Colab بسیار دشوار است.
- محدودیت در دیتاست‌ها:  
منابع داده‌ای موسیقی همراه با توضیح متنی کم هستند یا حق کپی‌رایت دارند.
- سازگاری نمونه‌برداری (Sampling Rate):  
مدل‌های مختلف (تولید و ارزیابی) ممکن است نرخ‌های متفاوتی داشته باشند.





## جمع‌بندی:

- رویکردهای پردازش زبان طبیعی برای تولید موسیقی پیشرفت زیادی کرده‌اند؛ از روش‌های ساده مبتنی بر قواعد تا مدل‌های پیچیده ژنراتور مانند MusicLM و MusicGen.

- ارزیابی موسیقی تولیدی نیازمند هم معیارهای آماری و هم شنیداری انسانی است و پیاده‌سازی موفق نیازمند درک دقیق از ارتباط میان زبان، احساس و ساختار موسیقی است.

## پیشنهادهای آینده:

1. ترکیب ارزیابی انسانی و ماشینی برای سنجش بهتر خروجی‌ها.
2. ساخت دیتاست‌های چندزبانه و سبک‌محور برای آموزش مدل‌های تطبیق‌پذیرتر.
3. بهینه‌سازی مدل‌ها برای محیط‌های کم‌منبع.
4. افزایش کنترل‌پذیری خروجی موسیقی با افزودن پارامترهای سبک، ریتم، ساز و ... در پرامپت.
5. استفاده از موسیقی‌درمانی یا آموزش موسیقی به‌عنوان کاربردهای عملی مدل‌های متن به موسیقی.
6. ...؟



1. Yoonjin Chung, Pilsun Eu, et 21 Feb 2025: “KAD: No More FAD!”
2. copet, J., et al. (2023). “Simple and Controllable Music Generation.” (MusicGen model card).
3. <https://huggingface.co/facebook/musicgen-small>
4. ...

با تشکر از توجه شما

[Aminkianiworkeng@gmail.com](mailto:Aminkianiworkeng@gmail.com)  
<https://github.com/M-Amin-Kiani/NLP-Proj>

