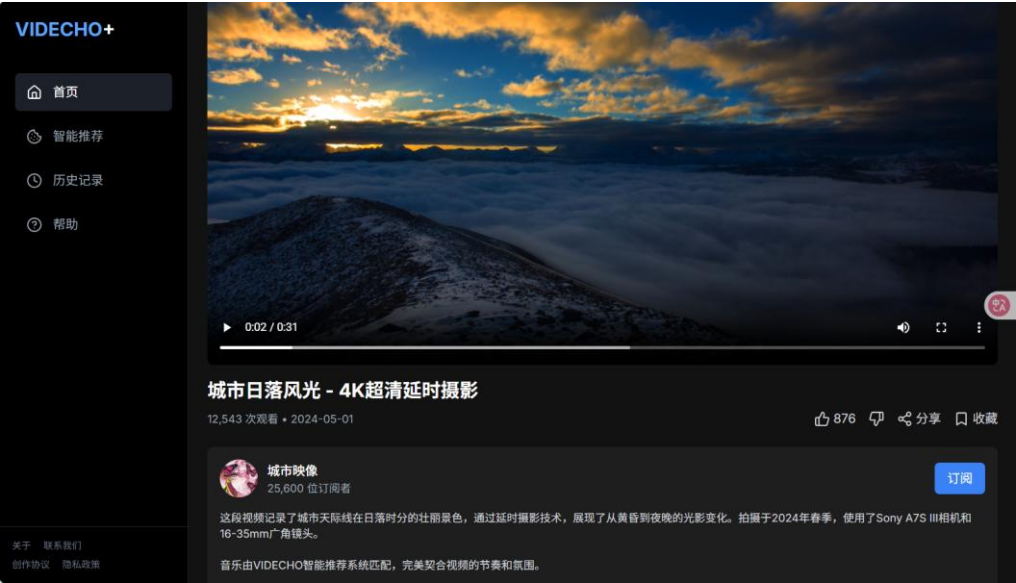
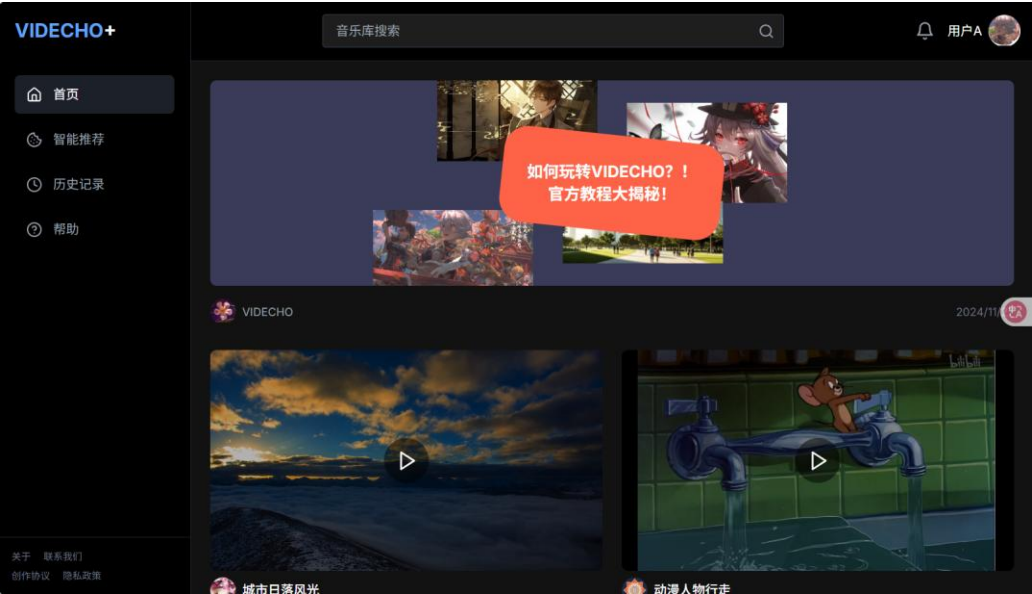
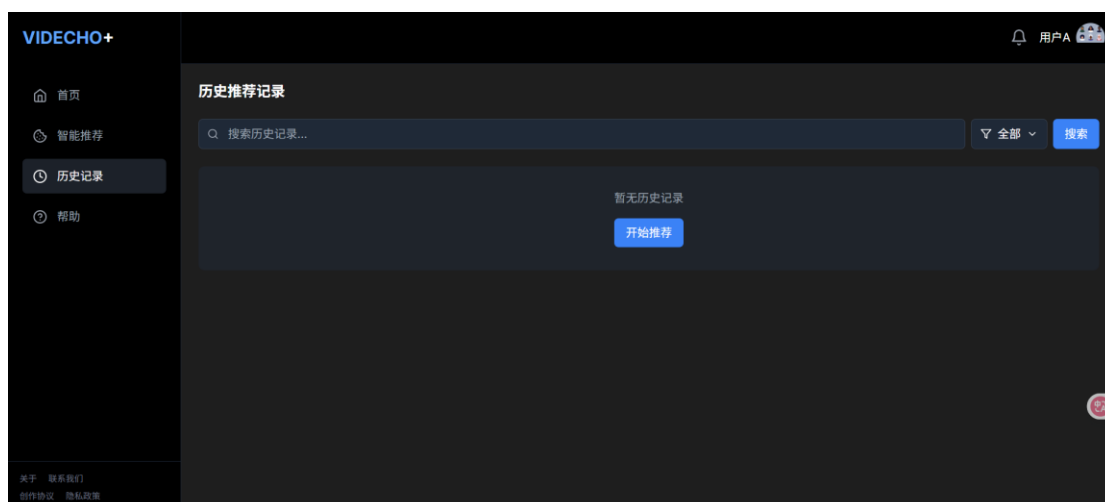
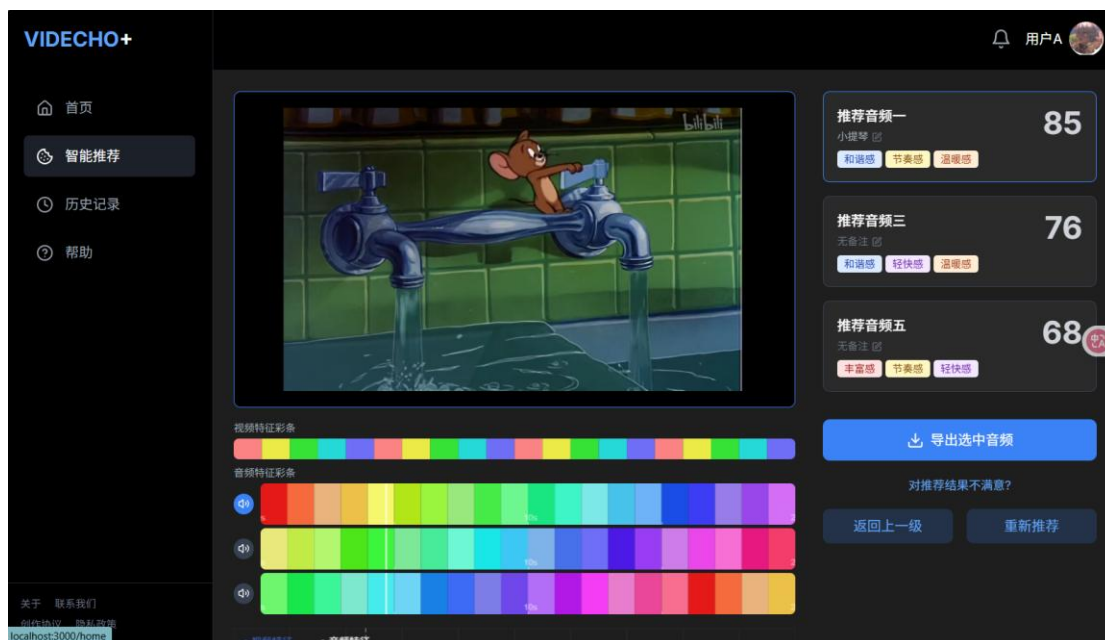
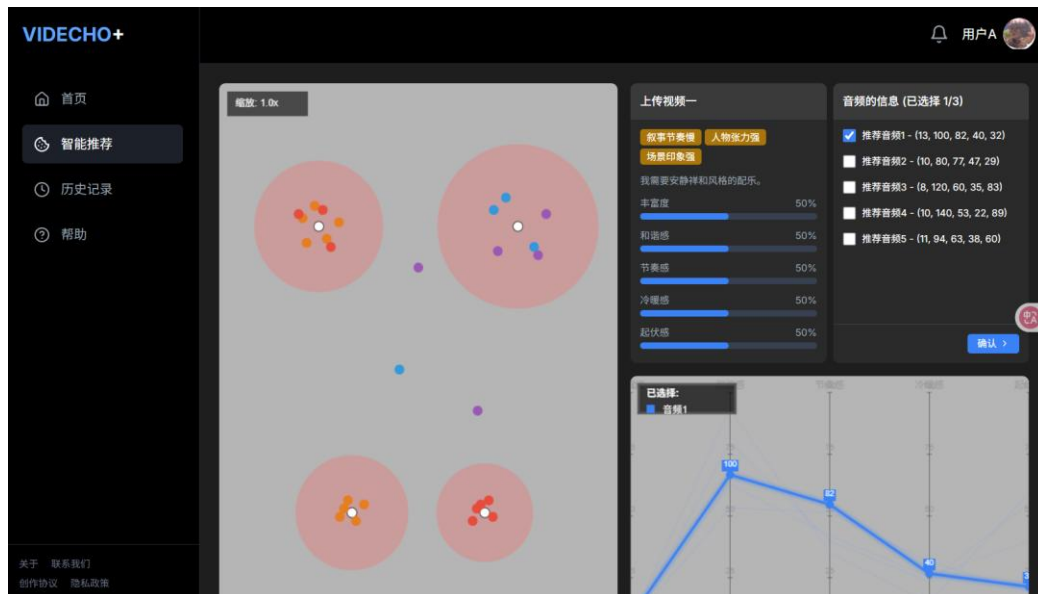
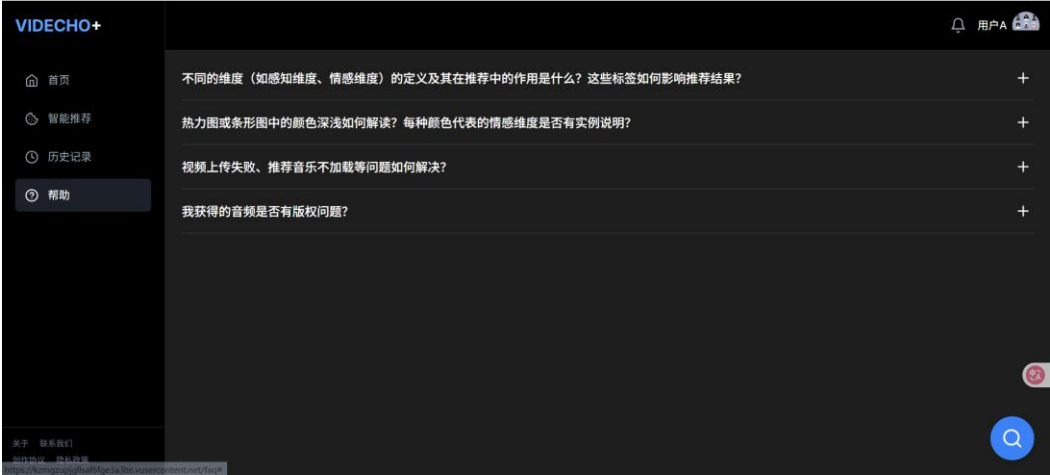
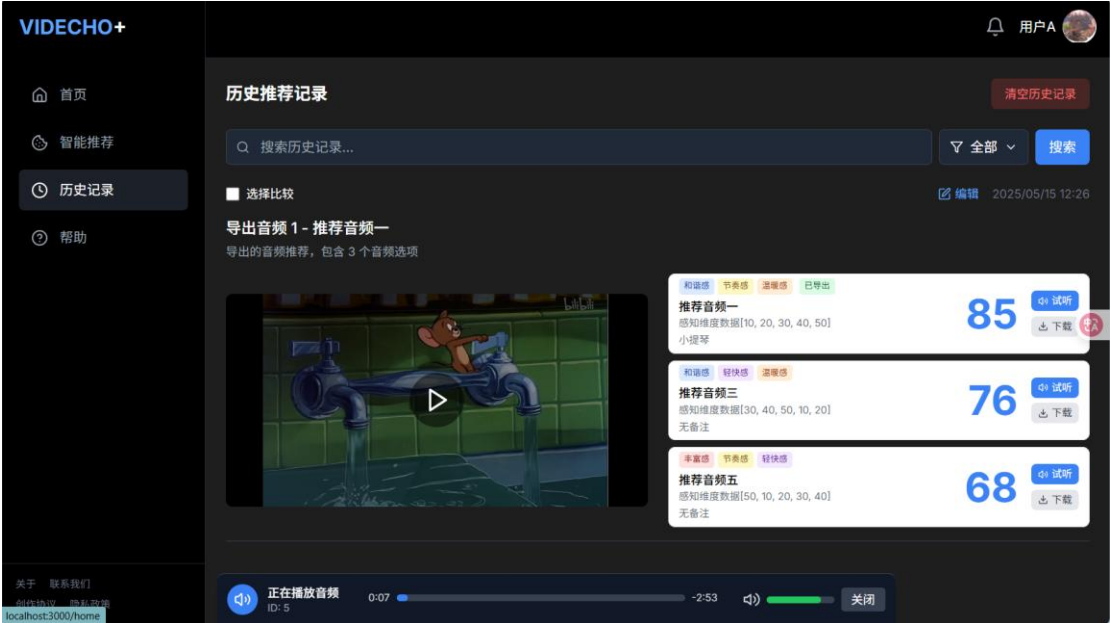
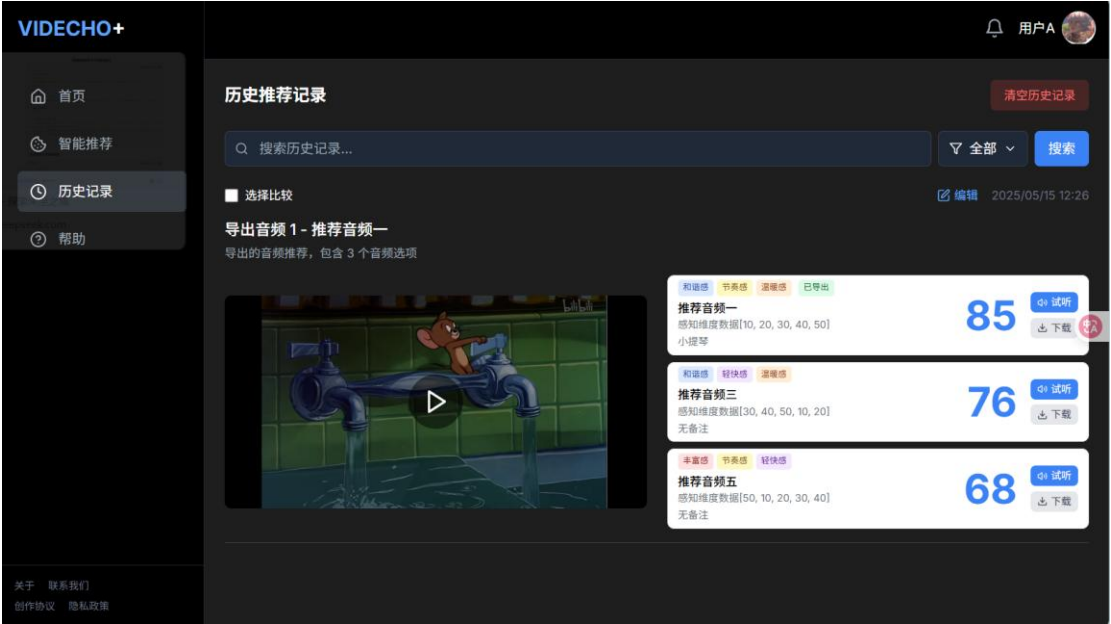
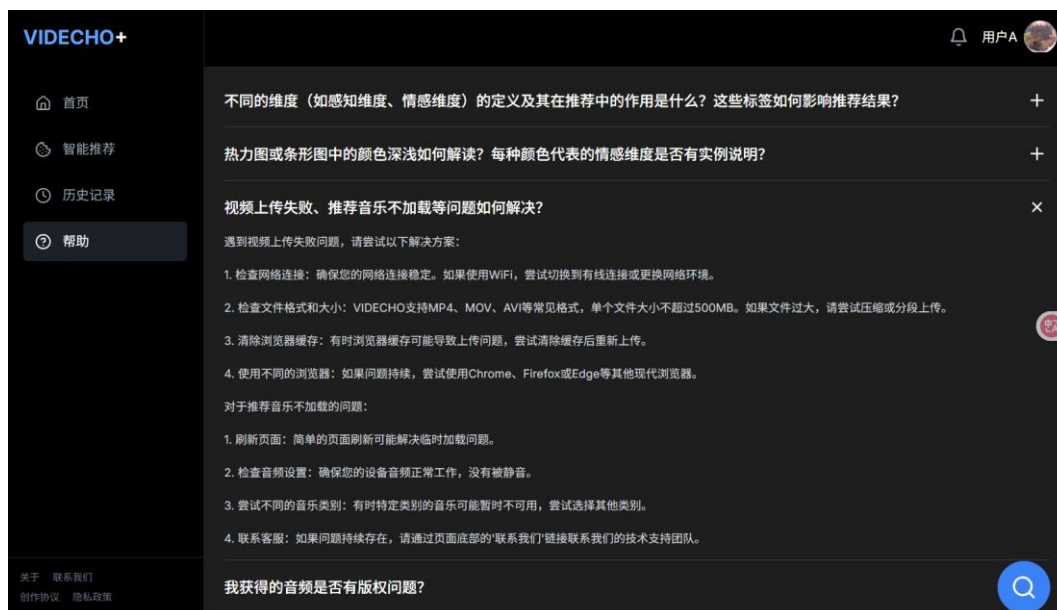


网页：<https://kzmgzupjgllsaf6fge3a.lite.vusercontent.net/home>









模型介绍：

1. 视频特征提取模型

目的：

提取视频中与配乐相关的语义和感知特征，作为后续推荐和匹配模型的输入。

输入数据：

原始视频片段（目前支持 MP4）

提取特征维度：

目前视频特征包含如下 11 个维度：

人物张力、场景张力、叙事节奏、叙事内容、形象亲缘、动作能量、色彩轻重、色彩冷暖、动作幅度、画面唤醒度、画面效价

方法与工具：

视觉感知特征：通过预训练的 CNN 提取帧级动作与色彩特征。

剪辑节奏特征：通过帧间变化速率、运动矢量分析等方式获得动作节奏感。

张力与情绪特征：结合面部表情、镜头运动与场景检测，使用轻量 Transformer 进行时间建模。

输出格式：

每个视频最终编码为一个十一维的向量，作为推荐模型的输入。

(下图为特征预测值的 R 方)

	R2
人物张力	0.735714
场景张力	0.688761
叙事节奏	0.837422
叙事内容	0.829261
形象亲缘	0.776947
动作能量	0.751509
色彩轻重	0.467869
色彩冷暖	0.764292
动作幅度	0.410173
画面唤醒度	0.814311
画面效价	0.752924

2. 视频-音频推荐模型

目的：

根据视频特征向量，预测其匹配的音频特征向量，用于推荐最合适的配乐片段。

模型结构：

主干网络：多层全连接神经网络（MLP），将视频特征映射到音频空间。

多任务学习设计：

主任务：音频特征回归（9 维-8 维）

辅助任务：音频情绪感知回归（视频 VA-音频 VA）

损失函数：

回归损失：MSELoss，用于保持预测音频向量在向量空间中的位置与真实音频接近。

分类损失：CrossEntropyLoss，提升模型对音频情绪属性的理解。

总损失：主任务和辅助任务加权组合，权重通过超参数搜索确定最优并进行使用。

训练方法：

数据对构成：每个训练样本由一组 (video_feat, audio_feat, audio_label) 构成。

优化器：Adam / AdamW，带有学习率衰减机制。

超参数优化：采用 Optuna 框架对 hidden_dim、dropout、learning_rate、任务权重等进行调参，优化目标为 Recall@3。

批次训练：支持 GPU 多线程 DataLoader，mini-batch 大小一般为 32~64。

输出：

每个视频将被预测为一个 10 维的音频特征向量，用于匹配最接近的真实音频。

```
输入的视频特征：
人物张力：60.7500
场景张力：66.2000
叙事节奏：54.0000
叙事内容：56.8000
形象亲缘：-16.5000
动作能量：66.8000
色彩轻重：0.3800
色彩冷暖：0.9700
动作幅度：0.2900
画面唤醒度：67.1400
画面效价：-29.2900

预测的音频特征：
Richness: -0.8372
Variation: 2.8622
Rhythm: -3.8342
Warmth: -5.1344
Harmony: -4.2335
brightness_score: 22.4444
Melody Trend Score: 50.9279
Rhythm speed score: 76.1797
音乐唤醒度：56.2267
音乐效价：-20.4412
```

3. 音频匹配模型

目的：

将推荐出的音频特征向量与音频数据库中的所有音频进行相似度计算，返回 Top-K 最相似的配乐结果。

输入数据：

数据库：已提取的音频特征集 audio_features.csv，每个音频为一个 10 维向量

预测向量：由推荐模型输出的带权重的 10 维音频向量（用于调整权重匹配）

匹配算法：

支持加权欧氏距离和余弦距离两种相似度计算方式。用户可以切换匹配方式。

支持特征维度选择：如只使用唤醒度和效价（VA）进行匹配。

支持用户权重自定义：用户可对各维度设置重要性权重，影响最终匹配得分。

默认使用欧氏距离进行匹配。

匹配过程：

1. 载入数据库中所有音频特征向量
2. 对每个样本计算与预测向量的距离
3. 排序取前 Top-5 个最相似音频供用户筛选

输出：

返回最相似音频的名称、索引、匹配距离等结果，可用于配乐推荐界面展示。

```
PS C:\Users\72416\Desktop\匹配> python match_audio.py --audio_csv audio_features.csv --pred_feat="-5.5244,-3.5814,-8.4680,-0.1249,-1.5983,20.7980,52.0154,77.2815,42.1747,-13.7280" --top_k 3
使用最后两列特征 = [ 42.1747 -13.728 ]
距离度量: Euclidean

Top-3 匹配结果:
1. 索引/名称: 1470 距离: 0.2874
2. 索引/名称: 48 距离: 1.1457
3. 索引/名称: 192 距离: 1.1841
```