

云智未来<sup>9th</sup>

第九届中国系统架构师大会  
SYSTEM ARCHITECT CONFERENCE CHINA 2017

# QQ音乐的个性化探索

腾讯音乐 / QQ音乐业务线 / 智能数据中心

 腾讯音乐娱乐 |  QQ音乐  全民K歌

# Overview



01 | 关于音乐，关于用户

02 | 音乐个性化的思考和演进

03 | 广告个性化的尝试

04 | AI时代一些好玩的尝试

# Overview



→ 01 | 关于音乐，关于用户

02 | 音乐个性化的思考和演进

03 | 广告个性化的尝试

04 | AI时代一些好玩的尝试



**QQ音乐** / 听我想听的歌

注册用户 **8亿**

DAU **1亿**



**全民K歌** / 你其实很会唱歌

注册用户 **4.6亿**

主力军 90后用户

Play on all your devices

iOS/Android

— PC

— H5

— IOS/ANDROID 设备

— 车载

— 智能音响

— TV



### 做最权威的正版音乐

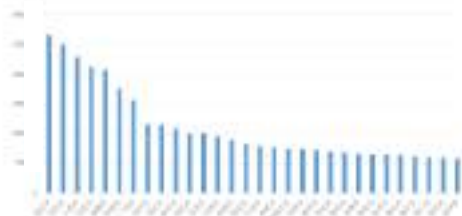
QQ音乐拥有的独家版权合作方包括华纳、索尼、环球、滚石、YG、LOEN、CUBE、福茂、新宝、华纳、华纳、少城时代、梦响当然等海内外优秀唱片公司20多家，累计达成版权战略合作方200多家，录制了超过1700万首的海量正版音乐，覆盖全曲之音。



## 一线城市阵地稳固，渠道下沉优势明显

QQ音乐平台 新增用户地域

QQ音乐用户市场排行TOP30



一线城市用户占比  
10.3%

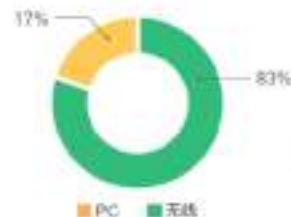
二线城市用户占比  
30.0%

三线及以下线城市用户占比  
59.7%

## 超级数字音乐航母在此腾飞

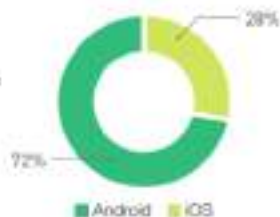
智能终端 | 网页+客户端 | 无线+PC | 安卓+IOS | WIN+MAC 全平台布局

平台播放分布



PC | 无线占比

移动端播放分布



iOS | Android占比





# Fantastic Baby

时间煮雨 我的歌声里 可惜没如果

滴答 泡沫告白气球 依然爱你

泡沫丑八怪李白 我好想你

青花瓷 幻听 小苹果 喜欢你

不再联系 平凡之路

一万个舍不得 七秒钟的记忆 闹够了没有

刚好遇见你 *Faded* 因为爱情

咱们结婚吧 我的好兄弟

# Always Online



时间维度：Aug.2017 – Oct.2017

## 巅峰音乐

-  《刚好遇见你》 李玉刚
-  《演员》 薛之谦
-  《小苹果》 筷子兄弟
-  《李白》 李荣浩
-  《平凡之路》 朴树
-  《丑八怪》 薛之谦
-  《告白气球》 周杰伦
-  《凉凉》 杨宗纬|张碧晨
-  《Faded》 Alan Walker
-  《默》 那英

## 巅峰专辑

-  《意外》 薛之谦
-  《不良少年》 徐良
-  《绅士》 薛之谦
-  《模特》 李荣浩
-  《三生三世十里桃花》 原声
-  《我很忙》 周杰伦
-  《魔杰座》 周杰伦
-  《万有引力》 汪苏泷
-  《JJ陆》 林俊杰
-  《刚好遇见你》 李玉刚

## 巅峰艺人

-  周杰伦
-  薛之谦
-  陈奕迅
-  张杰
-  林俊杰
-  许嵩
-  G.E.M. 邓紫棋
-  BIGBANG
-  徐良
-  张学友



# Overview



01 | 关于音乐，关于用户

→ 02 | 音乐个性化的思考和演进


03 | 广告个性化的尝试

04 | AI时代一些好玩的尝试

长期  
用户口碑，品牌调性

VS

短期  
业务KPI，case by case

人均听歌15%的  
大热之选

神曲到底是什么？

见过，听过的那些事儿

好听的歌是什么？

算法军备竞赛  
(CF+ Rules) VS DL



## 2011 - 2012



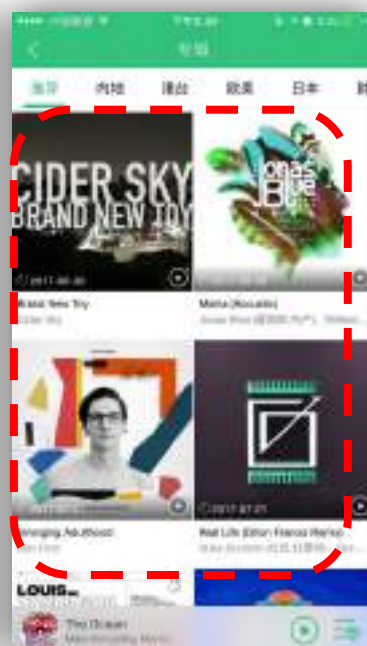
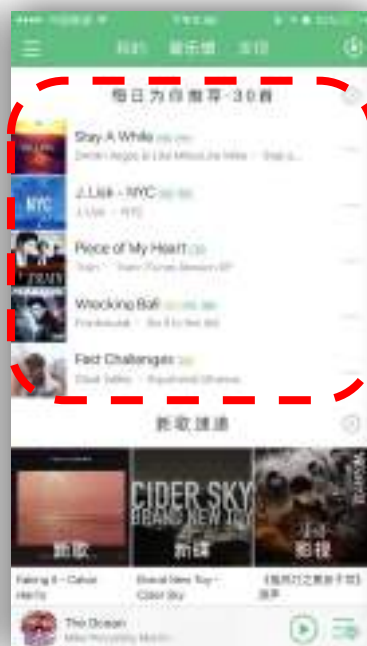
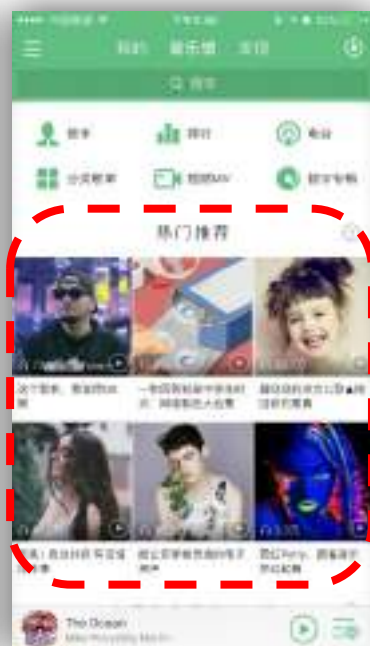
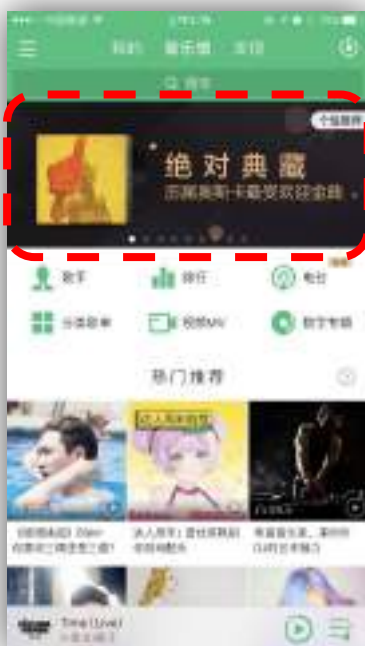
## 2013 - 2014

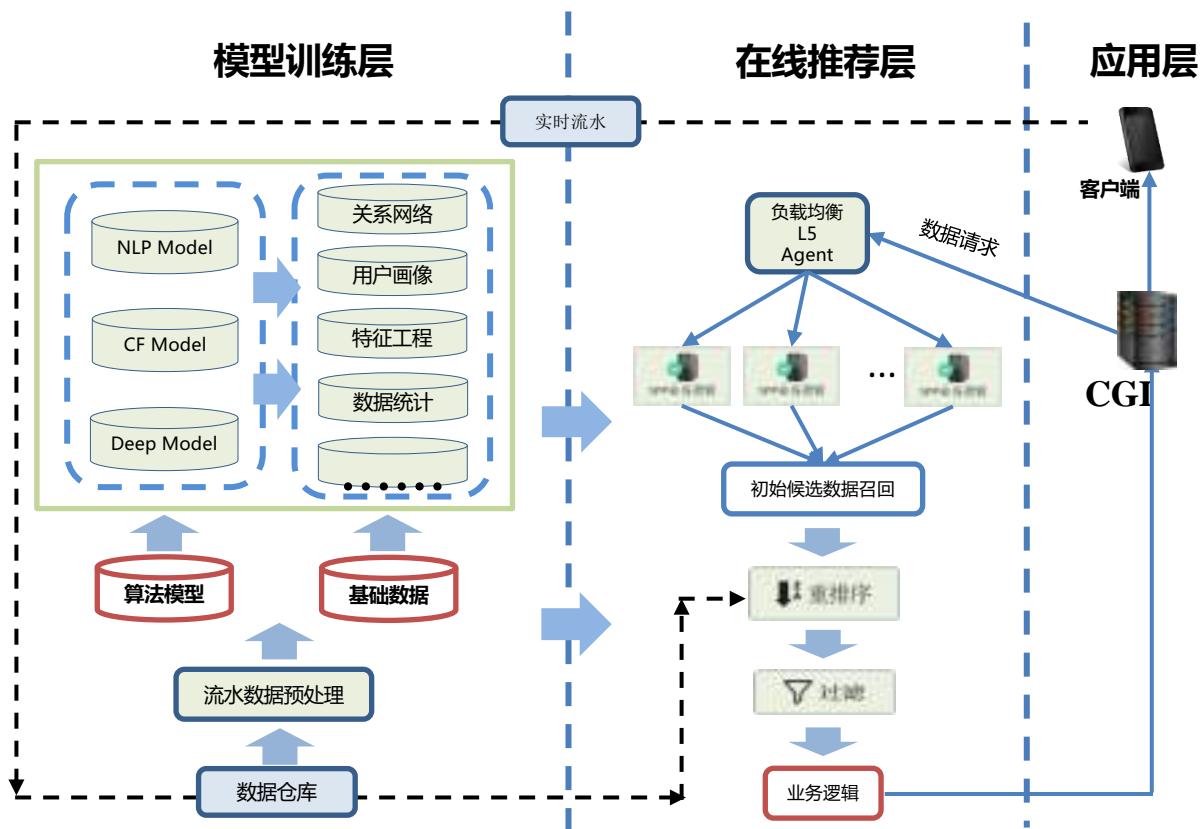


## 2014 - 2015



# 个性化推荐，依然在路上





## 中心化推荐

- 规则引擎
- 热门音乐
- 新音乐

## 引入个性化推荐

- 基于内容/标签推荐
- CF
- 基于用户长期兴趣歌手
- 基于用户长期兴趣流派

## 优化个性化推荐

- 在线实时架构
- 冷启动优化
- 用户特征工程优化
- 多目标推荐优化
- 内容特征优化 (歌单等文本类模型)

## AI相关探索

- 音乐大数据挖掘
- 深度神经网络
- 图像理解

## Pandora 专家人工标注



专业公司Gracenote以及学院派专业人员采用近2000种音乐元数据（流派、情感、主题标签）对每一首歌曲进行分析并标签化



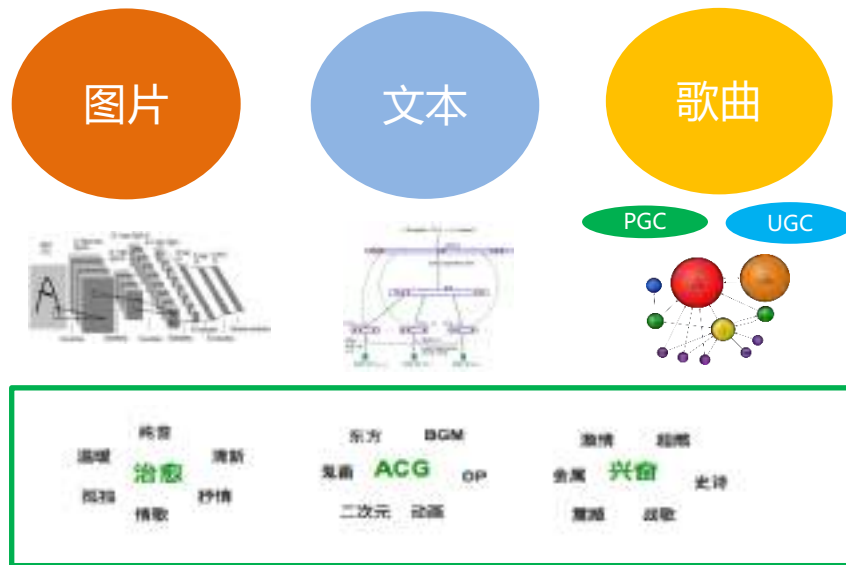
## 歌单知识体系



### ➤ 基础特征描述:

- 语种,流派,歌手,年代分布...
- 冷热程度
- 播放流水,收藏流水

### ➤ 语义特征描述:



## 歌单产品周期全面支撑

### 歌单投稿

日均机审占日审核量  
**70%+**

### 广场排序

歌单收听数量提升  
**56% -> 87%**

### 歌单推荐

音乐馆听歌显著增长

### 关联歌单

全面覆盖外部展示歌单





## 多平台联动

- QQ, 微信基础画像
- 腾讯视频
- 全民K歌
- 朋友圈/微博音乐分享数据等

## 海量用户行为数据挖掘

- DAU: 1亿+
- 单用户日均操作数据: 50亿+
- 每日歌曲播放: 十亿级

## 用户特征

- 用户基础信息: 性别、年龄、地域、学历...
- 音乐口味偏好: 歌手、流派、语言、年代
- 音乐行为偏好: 电台、收藏、下载、歌单、搜索
- 平台行为: 新增、留存、回流、活跃...

## 用户特征挖掘算法

- 听歌及操作流水, 时间衰减模型
- 自然语言处理: 对文本数据, 如评论、歌单的标题挖掘
- 噪声过滤: SPAM等
- 监督学习: 利用LR、GDBT等模型进行用户喜好预测

## 基础信息

- + 身份信息
- + QQ活跃信息(月)
- + QQ音乐活跃信息(月)
- + 异常信息

## 听歌习惯

- + 画像歌手
- + 语言能力
- + 流派能力
- + 心情偏好
- + 主题偏好
- + 节奏偏好
- + 听歌路径

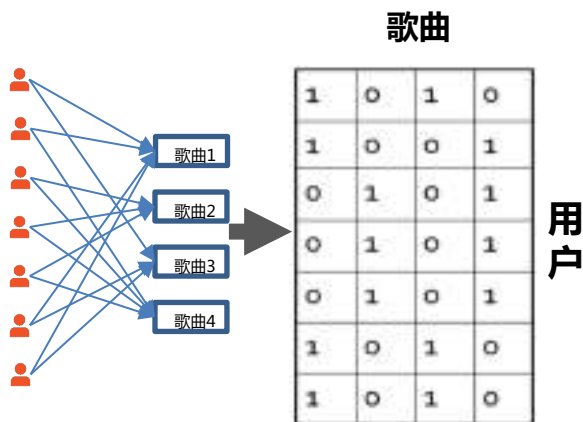
## 用户付费

- + 绿钻信息
- + 付费包(包绿)
- + 单曲、专辑、弹幕墙



## 浅层协同：

利用用户的听歌行为数据构建user-item矩阵，求取相似歌曲或相似用户，无须领域知识。

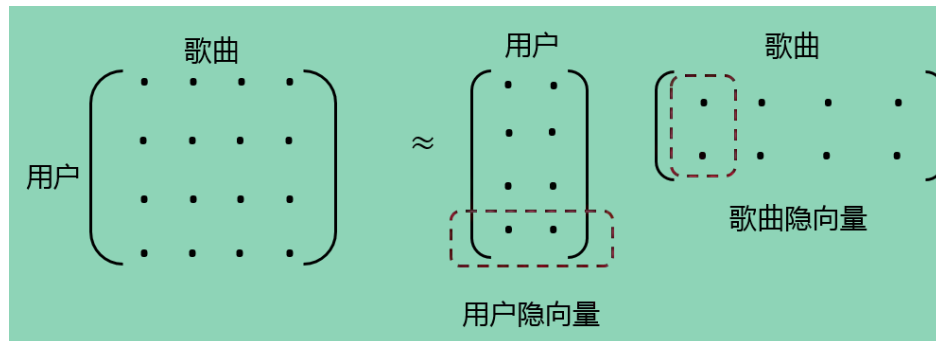


$$\text{sim}(i, j) = \frac{A \cap B}{A \cup B}$$

- 模型简单，准确度高，且可解析性好

## 隐因子模型：

Latent Factor Model：利用矩阵分解方法，求出用户和歌曲的隐特征向量。



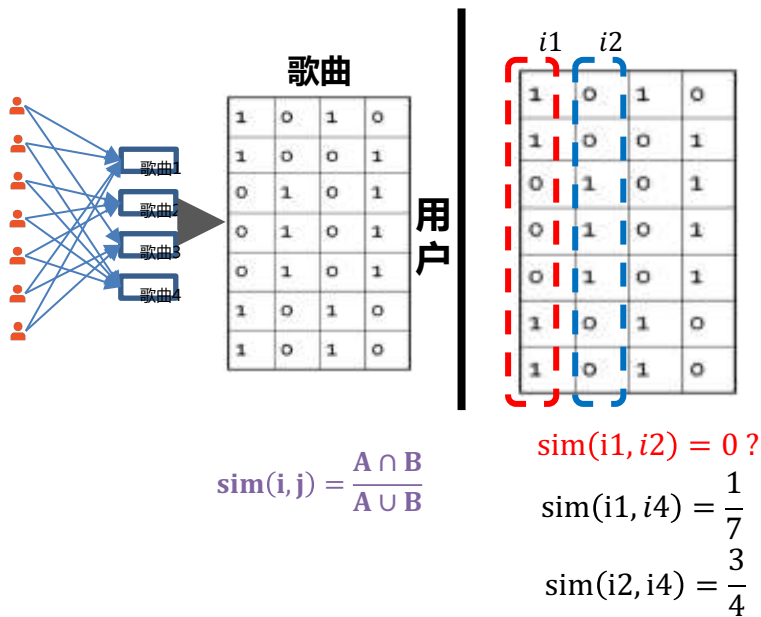
目标函数: 优化均方误差RMSE，使得预测评分与真实评分的均方误差最小

$$\min_{x^*, y^*} \sum_{u,i} c_{u,i} (p_{ui} - x_u^T y_i - \beta_u - \beta_i)^2 + \lambda \left( \sum_u \beta_u^2 + \sum_i \beta_i^2 \right)$$

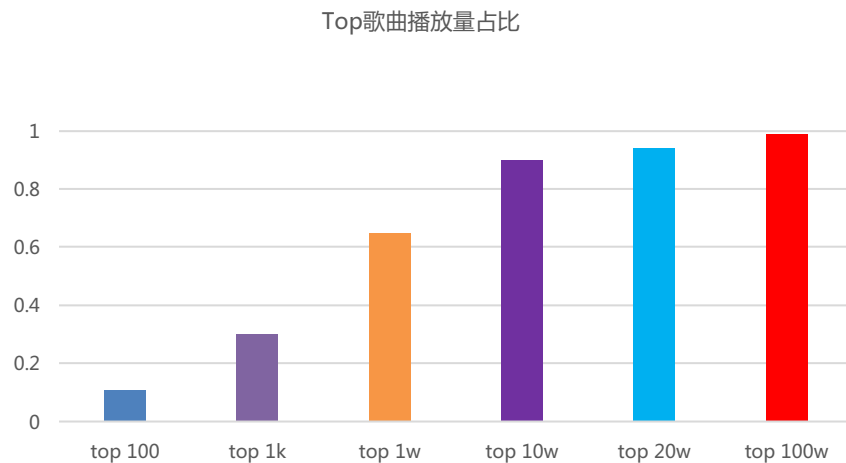
$p_{ui}$	一个二值变量，当 $r_{u,i} > 0$ 时，直接置为1，否则置为0；其中 $r_{u,i}$ 是用户 $u$ 对歌曲 $i$ 的得分
$\beta_u, \beta_i$	分别代表用户与歌曲的隐特征向量的偏移量
$x_u, y_i$	$x_u$ 代表用户的隐特征向量， $y_i$ 代表歌曲的隐特征向量
$c_{u,i}$	置信度，由于隐因子模型考虑的是隐式反馈，因此需要一个置信度来表示得分置信度 $c_{u,i} = 1 + \alpha * r_{u,i}$

## 协同模型(CF Model)的挑战

- 仅能发现浅层特征。
  - 缺点一：推荐的歌曲风格单一，缺乏新鲜感
  - 缺点二：只挖掘浅层的特征，无深层的特征



- 马太效应明显，Top100万歌曲占据了总收听量的90%+；
- 基于用户行为召回的数据，多以热门数据为主，如何跳出热歌圈子，挖掘长尾歌曲？
- 亿级用户的协同计算性能问题。



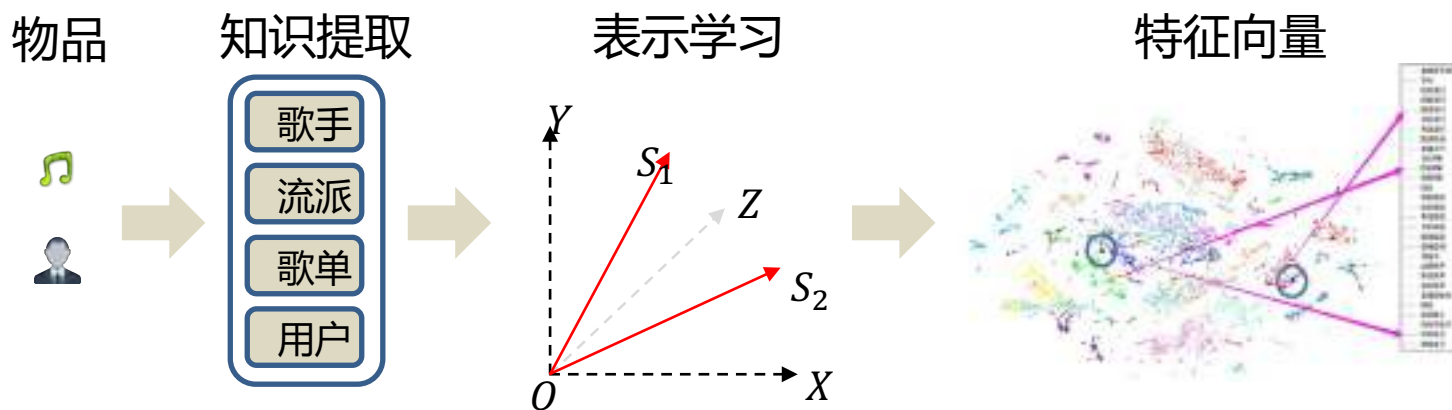
## • 相似度计算公式小优化

### - 引入IUF(Inverse User Frequency)

By John S.Brees, David Heckerman, Carl Kadie <Empirical Analysis of Predictive Algorithm for Collaborative Filtering>

$$sim(i, j) = \frac{\sum_{u \in N(i) \cap N(j)} \frac{1}{\log(1 + |N(u)|)}}{\sqrt{|N(i)| |N(j)|}}$$

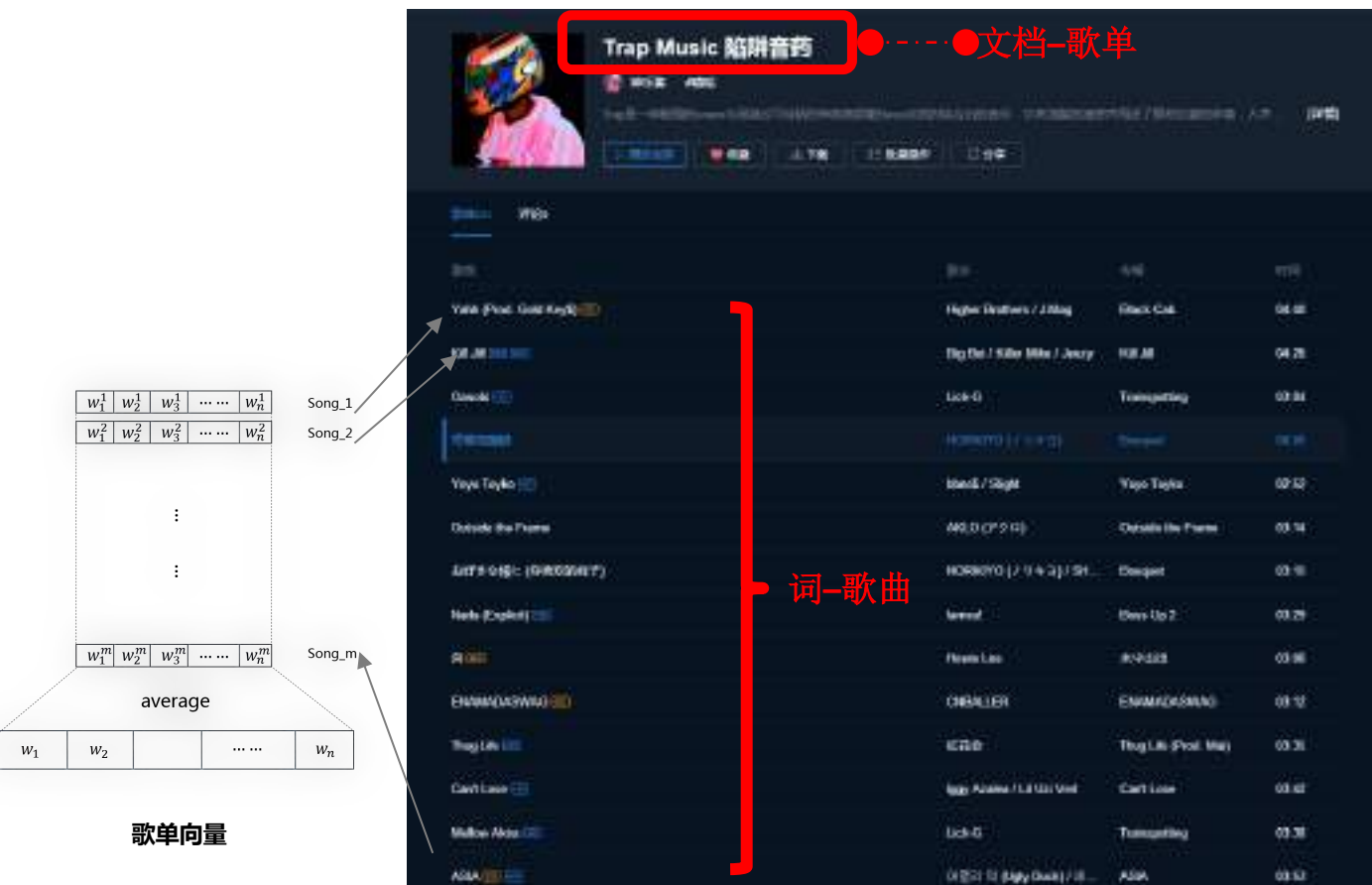
## • 离线计算度量优化



如何评价不同的特征向量表示的质量：

- 线上环境评测：相似单曲电台等实测场景。人均听歌时长，单曲听歌时长等。
- 模型的评测：特征向量的表示，应该能够使得在同一个流派下，歌曲之间的特征向量距离尽量接近。

$$\min \left( \frac{1}{m} \sum_v \|s_i - s_j\|^2 \right)$$



与CF Model相比，NLP Model的应用场景多种多样，既可以用于作为推荐数据召回建模，也可以用于特征提取建模

### 主动热度降权：

歌单数据中存在大量的小众歌曲，有利于进行长尾推荐；并且受到大盘听歌流水的影响较小，降低了噪音对模型的训练的影响

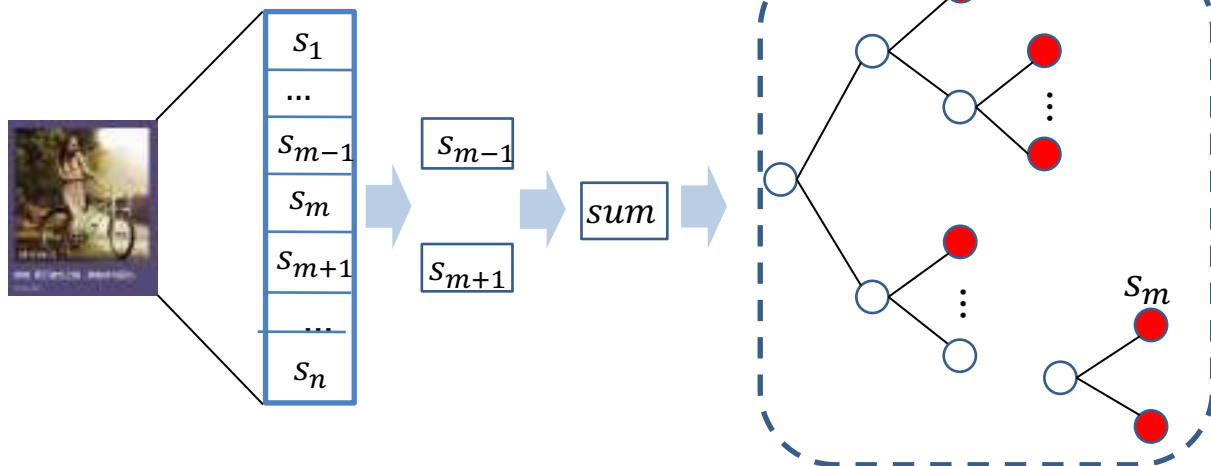
### NLP Model的主体思想：

将歌单作为文档，通过word2vec求取每一首歌曲单词的词向量表示

### 数据增强：

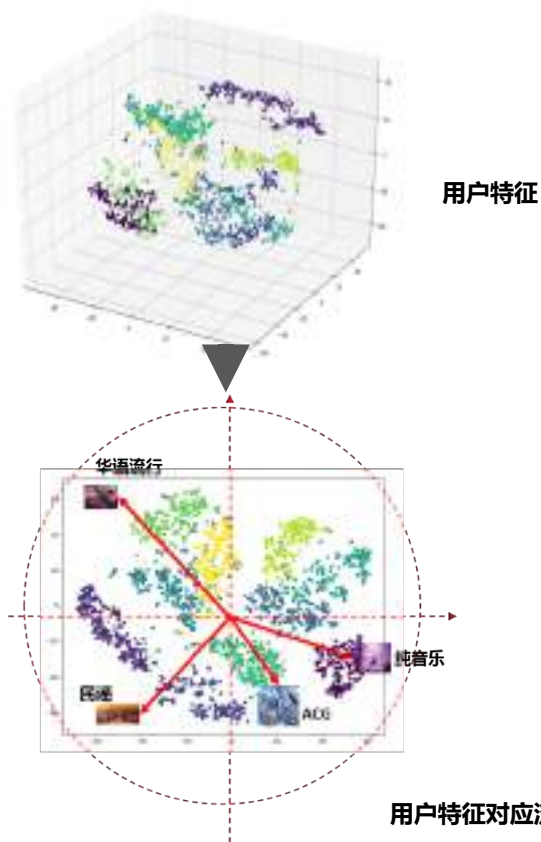
精选歌单质量好，但数量少，通过组合歌单来扩展语料库

Word2Vec



## Manifold Learning

利用t-SNE对高维的用户特征和歌曲特征进行降维:



- 将歌单作为文档，通过word2vec求取每一首歌曲的词向量表示，倾向于选择热门的歌曲做负样本进行训练 (negative sampling)。越热门的歌曲，离根节点越近。

用户没有点击某一首歌曲，通常有两种原因：

- 一是不知道有这首歌曲；
- 二是不喜欢这首歌曲

对于热门歌曲来说，显然第二种的可能性更高，这也是将热门物品作为负样本的合理性原因

## ● 一些阶段效果

### ➤ 产品指标的提升

- 听歌人数：Android，Iphone平台提升**20%+**
- 人均听歌时长：Android，Iphone平台提升**15%+**

### ➤ 一些用户好评



## Deep Neural Networks for YouTube Recommendations

Past Covington, Jay Adams, Emre Sargin  
Google  
Mountain View, CA  
{covington, ja, msargin}@google.com

### ABSTRACT

The following presentation of the biggest wide and deep neural network based recommendation system for YouTube, in this paper, we describe the system at a high level and focus on the Architecture, performance improvements through deep learning. The paper is split according to the classic recommender information retrieval (IR) model, so that a deep candidate generation model and then describe a separate deep ranking model. We also provide practical lessons and insights for real time deployment, training and monitoring a model on a distributed data center with systematic learning request.

### Keywords

recommendation system, deep learning, scalability

### 1. INTRODUCTION

YouTube is the world's largest platform for creating, sharing and discovering video content. YouTube recommendations have proven to be highly more than a billion more distinct personalized content from an overwhelming catalog of videos. In this paper we will focus on the massive, multi-deep learning based neural net on the YouTube video recommendation system. Figure 1 illustrates the recommendations on the YouTube mobile app home.

Recommending YouTube videos is extremely challenging due to the large scope of requests.

- **Scale:** Many users try recommendation algorithms on a daily basis, which results in a high volume of requests. It is important to ensure that the system can handle the high volume of requests and that the system is scalable for handling YouTube's massive user base and requests.
- **Feedback:** YouTube has a rich diverse range of video content, which is updated frequently. The recommendation system should be capable of handling the massive volume of new content, as well as the large volume of user requests.

YouTube mobile app is built on top of a set of APIs to recommend videos on a global scale. We provide a high-level overview of the architecture and the components of the system. We describe the architecture of the system, the data sources, and the components of the system. We describe the architecture of the system, the data sources, and the components of the system.



Figure 1. Recommendations displayed on YouTube mobile app home.

with well-understood IR model can be substituted from an exploration/exploitation perspective.

- **State:** Historical user behavior on YouTube is inherently difficult to predict due to sparsity and a number of stochastic external factors. We study different the general trends of user behavior and learn from the past to predict future behavior. Furthermore, context associated with content is particularly important and a well-defined ontology. The algorithm used to be robust to these factors by characterizing it as learning data.

In conjunction with other product teams across Google, YouTube has undergone a fundamental paradigm shift to search using deep learning as a general purpose solution for nearly all learning problems. Our system is built on Google Brain, a deep learning framework based on TensorFlow [1]. TensorFlow provides a flexible framework for implementing and training deep neural network architectures using hardware that is not tied to a specific vendor. Our models learn approximately one billion parameters and are trained on thousands of machines of examples.

To learn more about neural networks, see our blog post.

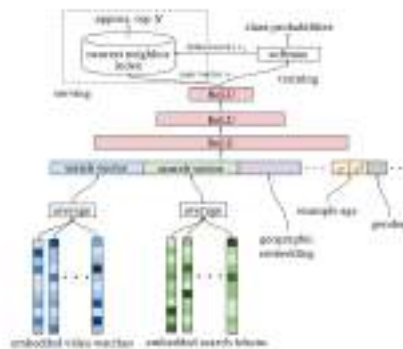


Figure 4. Deep candidate generation model architecture showing embedded sparse feature representation with video features. Embeddings are generated from combinations of candidate variable about tags of video files and their video content variables by input to the hidden layers. All hidden layers are fully connected. In training, a candidate loss is imposed with gradient descent on the output of the candidate softmax. At serving, an approximate nearest neighbor lookup is performed to generate thousands of candidate video recommendations.

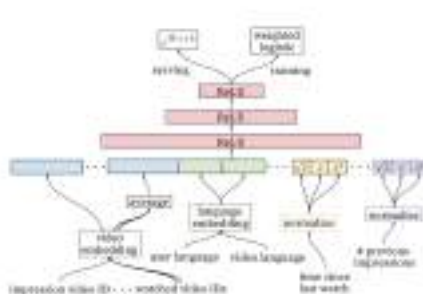
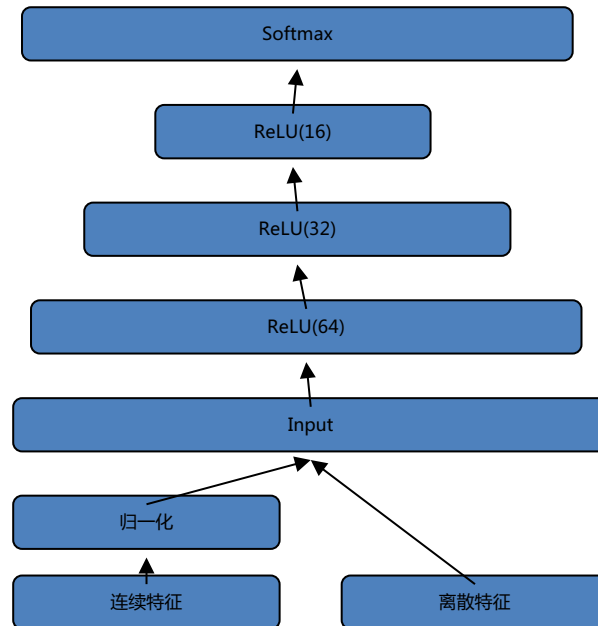


Figure 5. Deep ranking network architecture depicting embedded features from video-level and user-level with shared embeddings and generic of unimodal candidate features. All factors are fully connected. In practice, thousands of machines are fed into the network.





# Overview



01 | 关于音乐，关于用户

02 | 音乐个性化的思考和演进

→ 03 | 广告个性化的尝试

04 | AI时代一些好玩的尝试



## MusicBoss精准营销平台

- 产品运营自助配置广告；
- 效果追踪等一站式闭环管理；
- 定向投放与模型投放相结合；

用户包定向筛选

广告投放

活动模型

5000

1000

100

□冷启动：使用CF模型，针对特征标签少的用户及活动进行冷启动；

□特征拓展：利用word2vec等算法，进行特征维度拓展，并计算人群lookalike；

□活动推荐：排序模块使用Xgboost等模型，进行特征离散化及线上实时预测；

□线上优化：使用FTRL等算法，根据用户的反馈数据实时优化模型参数；

## 我们的在线广告探索：



### 生长阶段

2015年前

QQ音乐率先推出会员制，数字专辑等多种付费模式，推动音乐行业正版化，内部业务广告需求增长；同时逐步开放外部合作广告。**缺乏统筹，野蛮生长。**



### 产品化阶段

2015-2016年

逐步开始规范各个广告位的投放内容和形式。**纯人工运营，后台逐个需求开发。**



搭建广告统一管理平台

### 平台化阶段

2016-2017年

搭建了音乐广告管理平台，对所有广告位进行统筹管理，整合广告业务相关功能，引入推荐算法。**平台化整合各个功能，完成了推荐算法、数据分析、投放策略上的迭代升级。**



### 自动化阶段

2017年

加入了更多自动化功能，包括新增广告自动建模，流量自动分配等。**产品只需配置新广告，制定投放基本策略，平台在投放周期内会进行流量控制。**

## 平台化阶段 2016-2017年

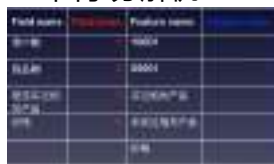
### APP内部广告

数字专辑  
付费音乐包  
外部合作  
会员推广

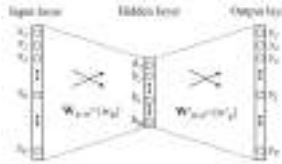


### 计算模块—核心算法变迁

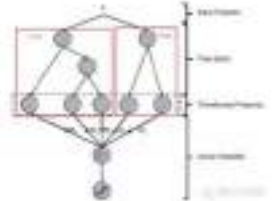
FM因子分解机



item2vec



GBDT+LR



XGBOOST



### 排期模块—运营与模型结合

定向投放

模型投放

定向+模型

频控投放

### 分析模块—全流程实时化

画像实时分析—Hermes (腾讯)

用户画像近10亿数据(约10TB)分析对比

	hadoop离线分析查询	实时多维分析
Count	15min+	1-3s
Sum	25min+	5-6s
Group by	26min+	1-3s

模型参数实时调整—FTRL



投放效果实时监控—TRC



## 广告模型自动建模和优化：

新上线广告，进行短期随机投放后，可**自动建模**；  
已上线广告，每隔一段时间会根据反馈数据，自动调整模型；  
自动化任务流如下：

### Assemble

特征集成，负责将样本和新特征进行集成

### Transform

特征转换，负责做特征常用转换，比如特征离散，特征交叉，tf-idf

### Criteria

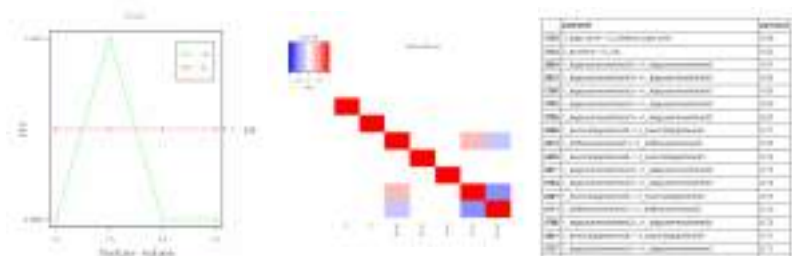
特征评估指标计算，包括entropy-ig, giniindex, entropy-igr, symmetry-uncertainty等

### Model

模型评估指标计算，包括auc, logloss, rmse等，以及输出特征全局重要度、树模型等



自动输出变量探索报告



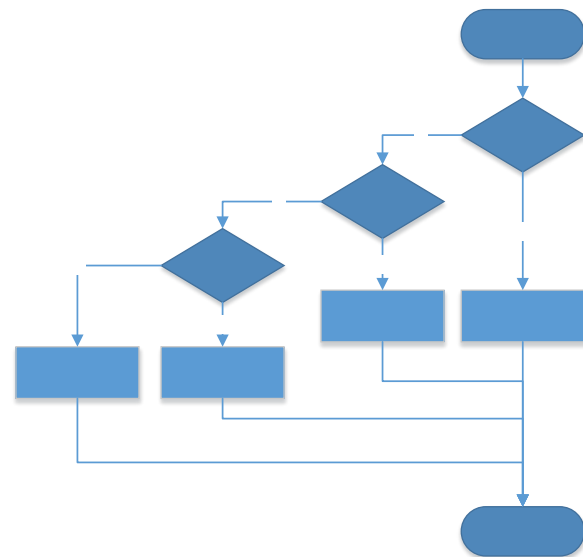
## 流量分发自动控制

广告id	广告名称	投放策略	所属平台	广告位置	投放模式	投放策略	投放平台	投放时间
1314	广告A	定向投放	QQ	手机桌面	定向投放	定向投放	iPhone/Android	2017-01-01
1315	广告B	定向投放	QQ	手机桌面	定向投放	定向投放	iPhone/Android	2017-01-01
1316	广告C	定向投放	QQ	手机桌面	定向投放	定向投放	iPhone/Android	2017-01-01
1317	广告D	定向投放	QQ	手机桌面	定向投放	定向投放	iPhone/Android	2017-01-01
1318	广告E	定向投放	QQ	手机桌面	定向投放	定向投放	Android	2017-01-01
1319	广告F	定向投放	QQ	手机桌面	定向投放	定向投放	Android	2017-01-01

为了处理人工运营、定向投放、智能推荐几种不同投放方式的矛盾，平台可进行不同投放方式的组合，例如：

- 1) 定向人群投放单一指定广告。
- 2) 定向人群进行多个广告的智能排序推荐；非定向人群按优先级投放。
- 3) 定向人群投放单一指定广告；非定向人群多个广告智能排序推荐。

当流量进入多个广告智能推荐分支时，又会根据模型效果，进行动态的流量调整。



# Overview



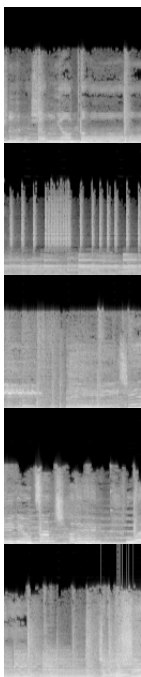
01 | 关于音乐，关于用户

02 | 音乐个性化的思考和演进

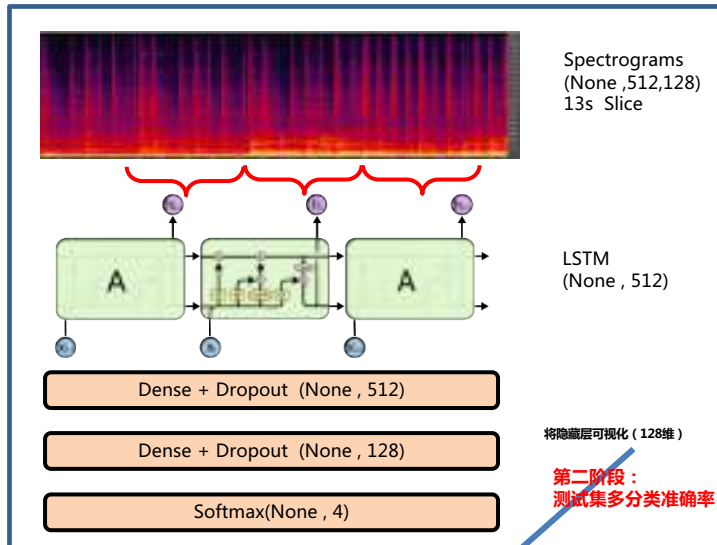
03 | 广告个性化的尝试

→ 04 | AI时代一些好玩的尝试

民谣  
电子  
爵士  
R & b  
摇滚

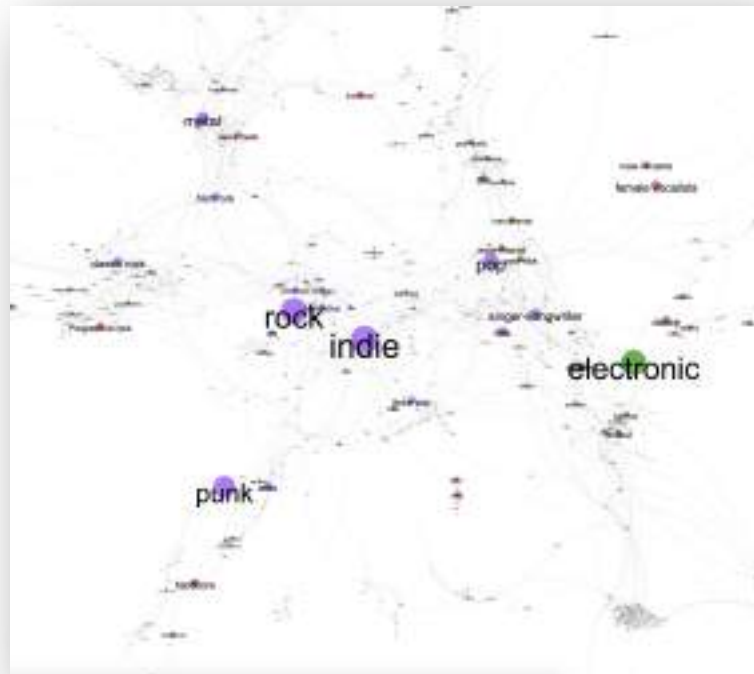
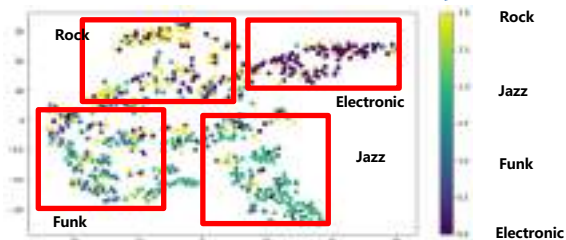


第一阶段：  
测试集多分类准确率~70%



将隐藏层可视化 (128维)

第二阶段：  
测试集多分类准确率~80%



By Echonetst

## 增量关联



## 统一的特征空间

主题：古典、乐器

实体：钢琴

抽象特征：[0.726, 0.032, 0.438, .....]

情绪：中性



## 存量分析（举例）

封面：



歌单：《纯音 | 50首轻缓闲适钢琴曲》

歌曲：《卡农》

相关实体：帕赫贝尔



# THANKS



腾讯音乐娱乐



QQ MUSIC  
DATA LAB