



# 基于 MCP 的网易云评论区 AI 理解系统

阅读时间：15-20分钟

作者：1mht

## 目录

### 1. 问题定义

- 1.1 场景背景
- 1.2 核心挑战
- 1.3 问题重新定义

### 2. 方案设计推导

- 2.1 用户真正要什么
- 2.2 技术约束分析
- 2.3 方案推导过程
- 2.4 核心设计原则

### 3. 系统架构

- 3.1 整体架构
- 3.2 采样系统详解
- 3.3 分析系统详解
- 3.4 数据库设计

### 4. 采样策略深度解析

- 4.1 为什么需要采样
- 4.2 采样结构设计
- 4.3 三级采样
- 4.4 动态分配策略
- 4.5 边缘情况处理
- 4.6 技术突破：cursor参数

### 5. 分层分析架构深度解析

- 5.1 为什么分层
- 5.2 Layer 0: 数据概览
- 5.3 Layer 1: 六维度信号

- 5.4 Layer 2: 验证样本
  - 5.5 Layer 3: 原始评论
  - 6. 六维度分析详解
    - 6.1 维度设计理念
    - 6.2 各维度详解
    - 6.3 置信度与局限性
    - 6.4 算法盲区识别
  - 7. 强制流程检查机制
    - 7.1 问题发现
    - 7.2 解决方案
    - 7.3 阈值设计
  - 8. 关键设计决策与权衡
    - 8.1 为什么不全量爬取
    - 8.2 为什么工具不下结论
    - 8.3 为什么需要对比样本
    - 8.4 为什么是100条阈值
  - 9. 典型使用流程
  - 10. 项目价值与反思
  - 11. Q&A
- 

# 一、问题定义

## 1.1 场景背景

网易云音乐评论区是中国互联网上一个独特的社会现象。它不仅仅是一个音乐讨论区，更是：

- **情感树洞**：用户在歌曲下分享个人故事、情感经历
- **文案博物馆**：大量精心撰写的"金句"被广泛传播
- **社交广场**：用户通过点赞、回复形成社区互动
- **时间胶囊**：评论跨越十几年，记录了不同时代的集体记忆

**我想做的事**：让 AI 能够"理解"一个评论区——不只是获取数据，而是真正理解听众的情感、话题、文化现象。

## 1.2 核心挑战

挑战	具体数字	为什么是问题
数据量巨大	热门歌曲 200万+ 评论	全量处理在时间和存储上都不现实
API 有限制	offset 最多访问 ~1100 条	官方接口只能获取最近评论，无法访问历史
时间敏感	用户等待容忍度 ~1分钟	全量爬取需要数小时，用户无法接受
Token 限制	200万条 $\approx$ 10亿 token	不可能把所有评论发给 AI
偏差风险	热评 $\neq$ 全部评论	只看热评会有平台推荐偏差

## 1.3 问题重新定义

**原始问题：**如何获取和分析200万条评论？

**重新定义：**如何用1000条评论让AI理解200万条评论的评论区？

这个重新定义非常关键——它把问题从"大数据处理"变成了"采样与推断"。

采样与推断的核心挑战是：

1. 如何保证样本有代表性？
2. 如何让 AI 知道样本的局限性？
3. 如何让 AI 在不确定性下做出合理判断？

# 二、方案设计推导

## 2.1 用户真正要什么？

当用户说"帮我分析这首歌的评论"时，他们真正想知道的是什么？

用户说：“帮我分析这首歌的评论”

- └─ 表面需求：获取评论数据
  - └─ 真实需求：理解评论区
    - └─ 情感层面
      - └─ 听众对这首歌有什么感受？
      - └─ 正面多还是负面多？
      - └─ 有没有什么特别的情感现象？
    - └─ 内容层面
      - └─ 他们在聊什么话题？
      - └─ 有没有反复出现的关键词？
      - └─ 有没有出圈的"金句"？
    - └─ 时间层面
      - └─ 评论区氛围随时间有变化吗？
      - └─ 是持续活跃还是有冷热周期？
      - └─ 有没有"复兴"现象？
    - └─ 文化层面
      - └─ 有什么有趣的文化现象？
      - └─ 评论区有什么独特的"氛围"？
      - └─ 和其他歌曲评论区有什么不同？

**核心洞察：**用户需要的是"理解"，不是"数据"。

## 2.2 技术约束分析

在设计解决方案之前，必须先理解技术边界：

约束类型	具体限制	影响
API offset 限制	超过 ~1100 条返回空	只能通过 offset 访问最近评论
API cursor 能力	可跳转任意时间戳	发现后成为突破口
请求频率	需要 0.5-1.5 秒间隔	大量请求耗时长

约束类型	具体限制	影响
热评数量	固定返回 15 条	平台精选，有偏差但有价值
AI Token	Claude 上下文有限	不能一次发送大量评论
用户耐心	~1分钟	超过就会放弃等待

## 2.3 方案推导过程

目标：让 AI 理解评论区

问题：AI 不能阅读 200万条评论

- |
- |— 方案 A：全量爬取后统计
  - | 优点：数据完整
  - | 缺点：耗时 2-6 小时，用户无法接受
  - | 结论：✗ 时效性不满足
- |
- |— 方案 B：只看热评（15条）
  - | 优点：最快，秒级响应
  - | 缺点：只代表平台推荐，严重偏差
  - | 结论：✗ 代表性不满足
- |
- |— 方案 C：随机采样（offset翻页）
  - | 优点：一定程度随机
  - | 缺点：offset限制，只能采最近评论，无历史覆盖
  - | 结论：✗ 时间覆盖不满足
- |
- |— 方案 D：分层采样（热评 + 最新 + 历史cursor跳转）
  - | 优点：快速（~1分钟）+ 时间覆盖完整 + 三种视角
  - | 缺点：仍是采样，必然有偏差
  - | 结论：✓ 可行，但需要透明告知偏差

选定方案 D 后的下一个问题：

采样必然有偏差，怎么办？

- |
- |— 方案 1：假装没有偏差，直接输出结论
  - | 风险：误导 AI，AI 过度自信
  - | 结论：✗ 不诚实
- |
- |— 方案 2：透明告知偏差，让 AI 自己判断
  - | 做法：告诉 AI 覆盖率、时间范围、算法局限
  - | 好处：AI 知道自己"不知道什么"
  - | 结论：✓ 诚实且有效

# 2.4 核心设计原则

经过推导，确立以下设计原则：

## 原则一：白盒设计

完全透明，让 AI 知道我们做了什么、没做什么、以及局限在哪里。

透明告知的内容：

- └─ 数据边界
  - | └─ 数据库有多少条评论
  - | └─ API 说总共有多少条
  - | └─ 覆盖率是多少
- |
- └─ 时间范围
  - | └─ 最早评论是哪年
  - | └─ 最新评论是哪年
  - | └─ 哪些年份有采样
- |
- └─ 算法局限
  - | └─ 情感分析的置信度
  - | └─ 已知的误判类型
  - | └─ 建议 AI 如何验证
- |
- └─ 采样方法
  - | └─ 热评来源
  - | └─ 最新评论来源
  - | └─ 历史评论来源

## 原则二：职责边界

明确 MCP 工具做什么、不做什么：

MCP 工具的职责：

✓ 做：

- 采样：从 200万 → 1000条
- 量化：计算情感分数、提取关键词、统计分布
- 透明告知：覆盖率、时间范围、算法局限
- 提供样本：让 AI 自己阅读原文

X 不做：

- 生成结论："听众很喜欢这首歌"
- 价值判断："这是一首好歌"
- 因果解释："因为疫情导致情感下降"
- 隐藏问题：假装数据完整

为什么工具不下结论？

因为我们使用的量化算法（如 SnowNLP 情感分析）有系统性偏差。SnowNLP 训练于商品评论，对音乐评论的"感伤式金句"经常误判为负面：

评论原文	算法判断	实际含义
"我失去了最好的朋友"	负面	深情怀念
"我恨这首歌让我哭"	负面	高度认可
"每次听都想起她，眼泪流下来"	负面	情感共鸣
"单曲循环到心碎"	负面	沉浸式喜爱

如果工具直接输出"负面情感占40%"这样的结论，会误导 AI。

正确做法：提供数据 + 提供样本 + 告知局限，让 AI 自己阅读、自己判断。

原则三：渐进式加载

AI 不需要一次性获取所有数据，而是按需加载：



Layer 0: 先看边界 → 决定是否需要采样  
Layer 1: 看量化信号 → 决定哪些需要验证  
Layer 2: 看验证样本 → 验证信号是否准确  
Layer 3: 看原始评论 → 深入特定问题

好处:

- 节省 Token
  - AI 可以在任何一层停下来
  - 避免信息过载
-

# 三、系统架构

## 3.1 整体架构

用户请求  
"帮我分析《晴天》的评论"



搜索 → 确认 → 入库  
search\_songs → confirm\_song\_selection → add\_song



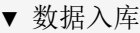
采样系统

职责：从 API 获取评论，存入数据库

工具：sample\_comments\_tool (quick / standard / deep)

固定结构：热评(15) + 最新(offset) + 历史(cursor)

动态分配：根据歌曲年龄调整比例



分析系统

职责：从数据库读取评论，渐进式分析

工具：Layer 0 → Layer 1 → Layer 2 → Layer 3

Layer 0：数据概览（边界信息）

Layer 1：六维度信号（量化指标）

Layer 2：验证样本（锚点+对比样本）

Layer 3：原始评论（按需筛选）



AI 生成报告  
(带证据引用，承认局限)

## 3.2 采样系统详解

### 为什么单独设计采样系统？

采样和分析是两个不同的关注点：

- 采样：一次性（或增量），关注如何从 API 获取数据
- 分析：可重复，关注如何从已有数据提取洞察

分开后的好处：

1. 避免重复请求 API
2. 可以增量采样
3. 分析可以反复运行

### 采样系统的输入输出：

输入：

- song\_id: 歌曲ID
- level: quick(200) / standard(600) / deep(1000)

处理：

1. 获取歌曲发行年份，计算年份跨度
2. 根据跨度选择分配策略
3. 采集热评（15条）
4. 采集最新评论（offset翻页）
5. 采集历史评论（cursor跳转每年）
6. 去重后存入数据库

输出：

- status: success/error
- actual: 实际采集数量
- samples: {hot: X, recent: Y, yearly: Z}
- coverage: {years\_span: N, years\_sampled: M}

### 3.3 分析系统详解

四层结构的设计理念：

## Layer 0: 边界层 — "我有什么数据?"

AI 需要知道:

- └─ 数据库里有多少条评论
- └─ API 说总共有多少条
- └─ 覆盖率是多少
- └─ 时间跨度是多少

AI 可以决策:

- └─ 数据是否足够分析?
- └─ 是否需要先采样?



## Layer 1: 信号层 — "数据说了什么?"

六个维度的量化信号:

- └─ sentiment: 情感分布
- └─ content: 主题关键词
- └─ temporal: 时间趋势
- └─ structural: 长度分布
- └─ social: 点赞集中度
- └─ linguistic: 评论类型

每个维度包含:

- └─ 量化数据
- └─ 置信度
- └─ 算法局限说明

AI 可以决策:

- └─ 哪些信号有价值?
- └─ 哪些信号需要验证?
- └─ 是否需要补充采样?



## Layer 2: 验证层 — "信号可信吗?"

两类样本:

- └─ 锚点样本 (不依赖算法)
  - └─ most\_liked: 最高赞

- | | | — earliest: 最早
- | | | — latest: 最新
- | | | — longest: 最长
- | |
- | | — 对比样本（发现算法盲区）
- | | | — high\_likes\_low\_score: 高赞但算法判低分
- | | | — low\_likes\_but\_long: 低赞但长文
- |
- | AI 可以做:
- | | — 阅读真实评论内容
- | | — 验证 Layer 1 的信号
- | | — 发现算法误判



Layer 3: 深挖层 — "需要看更多吗?"

- |
- | 按条件筛选原始评论:
- | | — year: 指定年份
- | | — min\_likes: 最低点赞数
- | | — limit: 返回条数
- |
- | 使用场景:
- | | — 验证某个年份的特殊现象
- | | — 深入某类高赞评论
- | | — 对比不同时期的评论风格

## 3.4 数据库设计

**设计理念:** 数据库 = 智能缓存

作用:

1. 避免重复请求 API
2. 支持增量采样（新评论追加，旧评论保留）
3. 记录元数据支撑透明度（覆盖率计算）

**核心表结构:**

```

-- 歌曲表
CREATE TABLE songs (
    id VARCHAR(64) PRIMARY KEY,           -- 歌曲ID
    name VARCHAR(255),                    -- 歌名
    artist VARCHAR(255),                  -- 艺术家
    album VARCHAR(255),                  -- 专辑
    publish_time BIGINT,                  -- 发行时间戳

    -- 缓存元数据
    cache_level VARCHAR(20),              -- none/basic/sampled
    cache_updated_at BIGINT,              -- 最后更新时间
    api_total_comments_snapshot INTEGER   -- API返回的评论总数快照
);

-- 评论表
CREATE TABLE comments (
    id INTEGER PRIMARY KEY AUTOINCREMENT,
    comment_id VARCHAR(64) UNIQUE,        -- 评论ID (API返回)
    song_id VARCHAR(64),                  -- 关联歌曲

    -- 评论内容
    content TEXT,                         -- 评论正文
    user_id VARCHAR(64),                  -- 用户ID
    user_name VARCHAR(255),               -- 用户名

    -- 评论元数据
    liked_count INTEGER,                  -- 点赞数
    timestamp BIGINT,                     -- 评论时间戳
    is_hot INTEGER DEFAULT 0,             -- 是否热评

    -- 管理字段
    is_deleted INTEGER DEFAULT 0,         -- 软删除标记
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);

```

---



## 四、采样策略深度解析

### 4.1 为什么需要采样

数学上的必然性：

假设：

- 热门歌曲有 200万 条评论
- 每条评论平均 50 字
- 总文本量 =  $200\text{万} \times 50 = 1\text{亿字} \approx 1.5\text{亿 token}$

Claude 的上下文限制：

- 即使是最大上下文也无法容纳
- 即使能容纳，处理成本也不可接受

结论：必须采样

采样的挑战：

挑战1：代表性

- 如何保证 1000 条能代表 200万？
- 如何避免系统性偏差？

挑战2：时间覆盖

- `offset` 只能访问最近评论
- 如何获取历史评论？

挑战3：透明度

- 采样必然有偏差
- 如何让 AI 知道偏差？

### 4.2 采样结构设计

三部分固定结构：

每次采样都包含：

1. 热评（固定15条）
来源：API hotComments
特点：平台算法精选，高社交价值
价值：代表"出圈"的内容，了解传播最广的观点
局限：有平台推荐偏差，不代表全部用户
2. 最新评论（offset翻页）
来源：offset=0 开始翻页
特点：当前时间点的评论
价值：捕捉当前评论区氛围，最新动态
局限：只能覆盖最近几个月
3. 历史评论（cursor跳转）
来源：cursor设为每年7月1日的时间戳
特点：跨年份分层采样
价值：保证时间线覆盖，发现长期趋势
局限：每年采样点有限

为什么是这三部分？

三个视角互补：

- 热评：社会共识视角
- 什么内容被最多人认可？
- 评论区的"代表作"是什么？
- 最新：即时状态视角
- 现在的评论区是什么氛围？
- 有没有新的趋势？
- 历史：时间演变视角
- 评论区氛围有变化吗？
- 不同年代的听众有什么不同？

## 4.3 三级采样

设计理念：只改数量，不改结构

级别	目标数量	适用场景	预计耗时
quick	200	快速预览、冷门歌曲	~30秒
standard	600	日常分析（推荐）	~60秒
deep	1000	深入研究、热门歌曲	~100秒

为什么是这三个数字？

quick (200):

- 热评15 + 最新100 + 历史85
- 适合快速了解，时间敏感场景
- 六维度分析勉强可用（部分维度置信度低）

standard (600):

- 热评15 + 最新175 + 历史410
- 六维度分析可靠
- 时间覆盖充分
- 平衡了质量和时间

deep (1000):

- 热评15 + 最新295 + 历史690
- 最高精度
- 适合需要深入研究的场景

## 4.4 动态分配策略

根据歌曲年龄调整比例：

歌曲类型	年份跨度	offset比例	历史比例	设计原因
新歌	≤2年	100%	0%	历史太短，全部用最新
中等新	3-5年	20%	80%	历史开始有意义

歌曲类型	年份跨度	offset比例	历史比例	设计原因
正常	6-10年	30%	70%	标准分配
老歌	>10年	30%	70%	限制采10年，避免太分散

分配示例 (standard, 8年跨度) :

```
target = 600
hot = 15 # 固定
remaining = 600 - 15 = 585

# 8年属于"正常"类型
offset_ratio = 0.3
yearly_ratio = 0.7

recent = 585 * 0.3 = 175 条 (offset采样)
yearly = 585 * 0.7 = 410 条
per_year = 410 / 8 = 51 条/年 (cursor采样)

最终分配:
  热评: 15
  最新: 175
  历史: 51 × 8 = 408
  总计: 598 (接近600)
```

## 4.5 边缘情况处理

```
# 冷门歌（评论总数 < 目标）
if api_total < target:
    # 全采，不分层
    return 全采（最多1000条）

# 新歌（≤2年）
if years_span <= 2:
    # 历史太短，全部用offset
    return 100% offset采样

# 老歌（>10年）
if years_span > 10:
    # 限制采10年，避免样本太分散
    effective_years = 10
    # 只采最近10年
```

## 4.6 技术突破：cursor参数

**问题：**网易云 API 的 offset 参数有限制，超过 ~1100 条返回空数据。

**发现：**通过 API 逆向分析，发现 weapi 接口的 cursor 参数是时间戳，可以跳转到任意历史时间点。

```
# cursor 参数使用方法
url = 'https://music.163.com/weapi/comment/resource/comments/get'

# 跳转到2020年7月1日
cursor = str(int(datetime(2020, 7, 1).timestamp() * 1000))

payload = {
    'rid': f'R_SO_4_{song_id}',
    'cursor': cursor, # 关键: 直接跳转到该时间点!
    'orderType': '1',
    'pageNo': '1',
    'pageSize': '20',
}

# 然后对payload进行weapi加密
encrypted = encrypt(payload)
response = requests.post(url, data=encrypted)
```

这个发现的价值：

对比	offset翻页	cursor跳转
可访问范围	最近~1100条	任意历史时间
时间覆盖	最近几个月	可达十几年
分层采样	不支持	支持
适用场景	最新评论	历史评论

## 五、分层分析架构深度解析

### 5.1 为什么分层

传统做法的问题：

传统做法：一次性返回所有分析结果

问题1：信息过载

- AI 收到大量数据，不知道关注什么
- 容易产生"确认偏差"—只看到想看的

问题2：Token浪费

- 很多信息 AI 可能不需要
- 全部发送浪费 Token

问题3：无法迭代

- 没有"发现问题→深入探究"的过程
- AI 无法主动决定下一步

## 分层的好处：

分层做法：渐进式加载，按需获取

好处1：AI 可以决策

- Layer 0 看边界，决定要不要继续
- Layer 1 看信号，决定要验证什么
- Layer 2 看样本，决定是否深入

好处2：节省 Token

- 如果 Layer 0 就发现数据不足，不需要继续
- 如果 Layer 1 没有异常，可能不需要 Layer 3

好处3：可追溯

- 每一层都有明确的输入输出
- AI 的判断过程可以追溯

## 5.2 Layer 0: 数据概览

**目的：**让 AI 知道数据边界，决定是否需要采样

**返回内容：**

```

{
  "layer": 0,
  "status": "success",

  "data_boundary": {
    "db_count": 587,           // 数据库有多少条
    "api_total": 324567,      // API说总共有多少条
    "coverage_rate": 0.18,    // 覆盖率
    "years_span": 8,          // 时间跨度（年）
    "earliest_year": 2016,    // 最早评论年份
    "latest_year": 2024       // 最新评论年份
  },

  "quality_assessment": {
    "sufficient": true,        // 数据是否足够分析
    "missing_years": [],       // 缺少哪些年份
    "recommendation": "数据充足，可以继续分析"
  },

  "ai_guidance": {
    "next_step": "get_analysis_signals_tool",
    "explanation": "数据量587条，覆盖8年，六维度分析可靠"
  }
}

```

## 强制检查:

```

if db_count < 100:
    return {
        "status": "must_sample_first",
        "blocking_reason": "数据量不足，当前仅 X 条，最低需要 100 条",
        "required_action": {
            "action": "sample_comments_tool",
            "params": {"song_id": song_id, "level": "standard"}
        }
    }
}

```



## 5.3 Layer 1: 六维度信号

**目的：**提供量化信号，让 AI 发现模式和异常

**六个维度：**

维度	分析内容	算法	置信度
sentiment	情感分布	SnowNLP	0.5-0.7
content	主题关键词	TF-IDF	0.7
temporal	时间趋势	年份统计	0.7
structural	长度分布	纯统计	0.9
social	点赞集中度	分布统计	0.85
linguistic	评论类型	规则分类	0.6

**返回结构示例：**

```
{
  "layer": 1,
  "dimensions": {
    "sentiment": {
      "distribution": {"positive": 0.67, "neutral": 0.18, "negative": 0.15},
      "confidence": 0.6,
      "limitation": "SnowNLP对感伤式金句可能误判为负面",
      "suggestion": "查看Layer 2的high_likes_low_score样本验证"
    },

    "content": {
      "keywords": ["青春", "回忆", "那年", "单曲循环", "眼泪"],
      "confidence": 0.7,
      "limitation": "TF-IDF只提取高频词, 无法理解语义"
    },

    "temporal": {
      "yearly_counts": {"2016": 45, "2017": 62, ..., "2024": 98},
      "trend": "increasing",
      "peak_year": 2024,
      "confidence": 0.7
    },

    "structural": {
      "length_distribution": {
        "micro(1-10)": 0.15,
        "short(11-30)": 0.35,
        "medium(31-80)": 0.30,
        "long(81-200)": 0.15,
        "very_long(200+)": 0.05
      },
      "confidence": 0.9
    },

    "social": {
      "gini_coefficient": 0.72,
      "top10_concentration": 0.43,
      "interpretation": "高度集中, TOP10评论占43%点赞",
      "confidence": 0.85
    }
  }
}
```

```
    },  
  
    "linguistic": {  
      "types": {"short": 0.35, "meme": 0.15, "story": 0.30, "review": 0.20},  
      "confidence": 0.6,  
      "limitation": "规则边界硬编码，分类可能不准"  
    }  
  },  
  
  "signals": [  
    {  
      "type": "high_likes_low_sentiment",  
      "description": "发现12条高赞但情感分数低的评论",  
      "action": "建议查看Layer 2的对比样本验证"  
    }  
  ]  
}
```

## 5.4 Layer 2: 验证样本

**目的：**提供真实评论样本，让 AI 验证 Layer 1 的信号

**两类样本：**

锚点样本（不依赖算法，纯粹基于客观指标）：

- └─ **most\_liked**: 最高赞评论
  - 用途：了解什么内容最受认可
  - 选择：按点赞数降序取前N条
- └─ **earliest**: 最早评论
  - 用途：了解早期听众的反应
  - 选择：按时间升序取前N条
- └─ **latest**: 最新评论
  - 用途：了解当前评论区状态
  - 选择：按时间降序取前N条
- └─ **longest**: 最长评论
  - 用途：了解深度讨论内容
  - 选择：按字数降序取前N条

对比样本（发现算法盲区）：

- └─ **high\_likes\_low\_score**: 高赞但算法判低分
  - 用途：发现情感分析误判
  - 选择：点赞>1000 且 情感分数<0.3
- └─ **low\_likes\_but\_long**: 低赞但长文
  - 用途：发现被埋没的深度内容
  - 选择：点赞<100 且 字数>100

**返回结构示例：**

```

{
  "layer": 2,
  "samples": {
    "anchors": {
      "most_liked": [
        {
          "content": "每次听到这首歌，都会想起那个夏天...",
          "likes": 125000,
          "timestamp": "2018-07-15",
          "sentiment_score": 0.85
        },
        ...
      ],
      "earliest": [...],
      "latest": [...],
      "longest": [...]
    },

    "contrast": {
      "high_likes_low_score": [
        {
          "content": "我失去了最好的朋友，这首歌是我们的回忆",
          "likes": 8500,
          "timestamp": "2020-03-21",
          "sentiment_score": 0.25,
          "analysis_note": "算法可能误判：内容是怀念，非负面"
        },
        ...
      ],
      "low_likes_but_long": [...]
    }
  },

  "ai_verification_tasks": [
    "验证high_likes_low_score样本是否为算法误判",
    "判断most_liked样本的共同主题",
    "对比earliest和latest评论的风格变化"
  ]
}

```

## 5.5 Layer 3: 原始评论

**目的：** 按需提供原始评论，支持深入探究

**使用场景：**

场景1：验证某年份的特殊现象

- 参数： `year=2020`, `limit=20`
- 用途：深入了解2020年的评论特点

场景2：查看更多高赞评论

- 参数： `min_likes=1000`, `limit=30`
- 用途：扩展`most_liked`样本

场景3：对比不同时期

- 调用两次： `year=2016` 和 `year=2024`
- 用途：对比早期和当前的评论风格

**返回结构：**

```
{
  "layer": 3,
  "filter": {
    "year": 2020,
    "min_likes": null,
    "limit": 20
  },
  "comments": [
    {
      "content": "...",
      "likes": 1500,
      "timestamp": "2020-05-12",
      "user_name": "****",
      "sentiment_score": 0.65
    },
    ...
  ],
  "total_matching": 89
}
```

## 六、六维度分析详解

### 6.1 维度设计理念

#### 为什么是这六个维度？

设计出发点：覆盖"理解评论区"的各个方面

**sentiment**（情感）

- 回答：听众对这首歌的情感反应是什么？
- 算法：SnowNLP

**content**（内容）

- 回答：他们在聊什么话题？
- 算法：TF-IDF + 规则

**temporal**（时间）

- 回答：评论区活跃度和氛围有变化吗？
- 算法：年份分组统计

**structural**（结构）

- 回答：评论的"形态"是什么？长评多还是短评多？
- 算法：纯统计

**social**（社交）

- 回答：点赞分布是什么？是精英控场还是大众狂欢？
- 算法：基尼系数、集中度

**linguistic**（语言）

- 回答：评论的"类型"是什么？玩梗/故事/乐评？
- 算法：规则分类

## 6.2 各维度详解

### sentiment (情感维度)

算法: SnowNLP

原理: 基于朴素贝叶斯, 训练于商品评论

输出: 每条评论一个0-1的分数 (0=负面, 1=正面)

统计方式:

positive:  $\text{score} > 0.6$  的比例

neutral:  $0.4 \leq \text{score} \leq 0.6$  的比例

negative:  $\text{score} < 0.4$  的比例

置信度: 0.5-0.7 (较低)

已知局限:

- 训练数据是商品评论, 音乐评论有不同特点
- "感伤式金句"被误判为负面
- 反讽、隐喻理解不准

应对策略:

- 提供置信度让 AI 谨慎对待
- 提供高赞低分对比样本
- 建议 AI 阅读原文验证



## content (内容维度)

算法: TF-IDF + 规则过滤

原理: 提取高频且有区分度的词

输出: TOP-N 关键词

处理流程:

1. jieba 分词
2. 去停用词
3. TF-IDF 计算
4. 规则过滤 (去掉"歌曲"、"好听"等通用词)
5. 返回 TOP-20

置信度: 0.7

已知局限:

- 只能提取高频词, 无法理解语义
- 可能遗漏低频但重要的内容
- 同义词会分散计数

应对策略:

- 提供原始样本让 AI 理解语义
- 关键词只作为参考, 不作为结论

## temporal (时间维度)

算法：年份分组统计

原理：按年份统计评论数量

输出：每年评论数、趋势判断

统计内容：

- 各年份评论数
- 增长/下降趋势
- 峰值年份
- 是否有"复兴"现象

置信度：0.7

已知局限：

- 只有数量趋势，没有内容变化
- 不知道"为什么"某年评论多
- 采样可能影响各年代表性

应对策略：

- 提供各年份样本让 AI 对比内容
- 趋势只作为参考

## structural (结构维度)

算法：纯统计

原理：统计评论长度分布

输出：各长度区间的比例

长度分类：

micro: 1-10字（表情、感叹）

short: 11-30字（简短感想）

medium: 31-80字（完整表达）

long: 81-200字（详细评论）

very\_long: 200+字（长文/小作文）

置信度：0.9（最高）

局限：几乎没有

应用：

- 长评多 → 评论区有深度讨论氛围
- 短评多 → 评论区更像即时反应
- 分布可以反映评论区"文化"

## social (社交维度)

算法: 分布统计

原理: 分析点赞分布的集中度

输出: 基尼系数、TOP-N集中度

指标:

`gini_coefficient`: 基尼系数 (0-1, 越高越不平等)

`top10_concentration`: TOP10评论占总点赞的比例

`top10_vs_median`: TOP10平均点赞 / 中位数

置信度: 0.85

解读:

基尼系数 > 0.7 → "精英控场", 少数评论获得大多数点赞

基尼系数 < 0.5 → "大众狂欢", 点赞分布相对平均

应用:

- 了解评论区的社交结构
- 判断是否有"意见领袖"

## linguistic (语言维度)

算法：规则分类

原理：基于长度、关键词、格式判断类型

输出：各类型的比例

类型：

**short**：短评（<15字，简单反应）

**meme**：玩梗/吐槽（≤30字，无故事特征）

**story**：故事/小作文（≥30字，有故事特征词）

**review**：乐评/鉴赏（≥15字，有音乐术语）

分类规则：

```
if len < 15: return "short"
if len <= 30 and no_story_feature and no_review_feature: return "meme"
if len >= 30 and has_story_feature: return "story"
if len >= 15 and has_review_feature: return "review"
return "short"
```

置信度：0.6（最低）

已知局限：

- 规则边界硬编码
- 复杂内容难以分类
- 可能有误分类

应对策略：

- 置信度最低，仅作参考
- 结合样本理解真实类型

## 6.3 置信度与局限性

为什么要标注置信度？

传统做法：返回数据，不说明可靠性

- AI 可能过度信任
- AI 不知道何时需要验证

我们的做法：每个维度标注置信度和局限

- AI 知道哪些可信、哪些需要验证
- AI 可以在报告中说明不确定性

## 置信度排序：

structural (0.9) > social (0.85) > content (0.7) = temporal (0.7) > linguistic (0.6) > sentiment (0.5)

## 6.4 算法盲区识别

**问题：**情感分析算法经常误判

**解决：**提供"对比样本"让 AI 发现盲区

对比样本：high\_likes\_low\_score

选择条件：点赞 > 1000 且 情感分数 < 0.3

这些评论的特点：

- 用户高度认可（高赞）
- 算法判断负面（低分）
- 说明算法可能误判

AI 阅读后可以：

- 判断是真的负面还是算法误判
  - 理解误判的模式（感伤≠负面）
  - 修正对 sentiment 维度的判断
-

# 七、强制流程检查机制

## 7.1 问题发现

早期测试的问题：

场景：用户说"帮我分析《晴天》"

AI 的实际行为：

1. 直接调用 Layer 1 查看信号
2. 用数据库里之前的旧数据（可能只有70条）
3. 输出分析报告

问题：

- 覆盖率极低（ $70/600000 = 0.01\%$ ）
- 时间覆盖不完整
- AI 不知道数据不足
- 用户看到的报告质量很差

根本原因：

AI 的倾向：

- 尽快给出答案
- 避免"麻烦"用户
- 不主动质疑数据质量

系统的问题：

- 没有强制检查数据充足性
- 允许 AI 跳过采样步骤

## 7.2 解决方案

在 Layer 0 和 Layer 1 设置阻断机制：

```
# Layer 0 强制检查
def get_analysis_overview(song_id):
    db_count = count_comments(song_id)

    if db_count < MIN_REQUIRED_FOR_ANALYSIS: # 100
        return {
            "status": "must_sample_first",
            "blocking_reason": f"数据量不足: 当前仅 {db_count} 条, 最低需要 100 条",
            "required_action": {
                "action": "sample_comments_tool",
                "params": {"song_id": song_id, "level": "standard"},
                "instruction": "必须先采样再继续分析"
            },
            "ai_instruction": "禁止继续调用 Layer 1/2/3, 必须先完成采样"
        }

    # 数据充足, 正常返回
    return {...}
```

三个检查点:

检查点	阈值	返回状态	行为
Layer 0	评论数 < 100	must_sample_first	阻断, 要求先采样
Layer 1	评论数 < 100	must_sample_first	阻断, 要求先采样
Layer 2	非 deep 级别	sampling_upgrade_prompt	提示是否需要更深入

7.3 阈值设计

为什么是100条?



分析各阈值的影响：

< 30 条：

- 统计分析基本无效
- 六维度信号不可靠
- 应该阻止分析

30-100 条：

- 基础分析勉强可行
- 部分维度置信度低
- 边界情况，可以允许但要警告

≥ 100 条：

- 六维度分析可正常运行
- 各维度置信度达到标注水平
- 可以输出可靠报告

结论：100 是平衡点

- 既保证基础分析可行
- 又不过于严格阻挡用户

## Layer 2 的升级提示：

```
# 非 deep 级别时，提示用户
if current_level != "deep":
    return {
        ...,
        "sampling_upgrade_prompt": {
            "should_ask_user": True,
            "current_level": current_level,
            "current_count": comment_count,
            "prompt_template": f"当前分析基于 {comment_count} 条评论（{current_level} 级别）。如需更精
```

# 八、关键设计决策与权衡

## 8.1 为什么不全量爬取

对比：

方案	时间成本	存储成本	用户体验	数据完整性
全量爬取	2-6小时	数GB	✗ 不可接受	✓ 100%
分层采样	30-100秒	几MB	✓ 可接受	~0.05%

决策：选择分层采样

理由：

1. 用户体验优先

- 用户不会等待数小时

- 即使愿意等，也不是好的产品体验
2. 边际效益递减

- 从0到1000条：信息量急剧增加

- 从1000到200万条：增加的信息有限

- 1000条足以覆盖主要模式
3. 透明度弥补

- 告诉 AI 覆盖率

- AI 可以在报告中说明局限

- 用户知道这是采样结果

## 8.2 为什么工具不下结论

传统做法：

传统分析工具的输出：

"情感分析结果：正面 67%，负面 15%，中性 18%"

"结论：听众普遍喜欢这首歌"

问题：

- 算法可能误判
- 结论可能误导 AI
- AI 无法质疑工具的判断

## 我们的做法：

我们的输出：

"情感分布：正面 67%，负面 15%，中性 18%"

"置信度：0.6"

"局限：SnowNLP 对感伤式金句可能误判为负面"

"建议：查看 high\_likes\_low\_score 样本验证"

"样本：[具体评论内容]"

好处：

- AI 知道置信度，不会过度信任
- AI 可以阅读样本自己判断
- AI 可以在报告中修正误判

## 具体案例：

评论："我失去了最好的朋友，这首歌是我们的回忆"

SnowNLP 判断：负面（分数 0.25）

原因：检测到"失去"、"回忆"等词

如果工具下结论：

- "这是一条负面评论"
- AI 信任工具，统计为负面
- 最终报告偏差

我们的做法：

- 返回分数 0.25，但标记为"高赞低分样本"
- AI 阅读原文，发现是深情怀念
- AI 在报告中纠正："虽然算法判断为负面，但阅读原文发现这是怀念式的正面情感"

## 8.3 为什么需要对比样本

**问题：**锚点样本（最高赞、最早等）都是基于客观指标，但无法发现算法盲区。

**解决：**增加对比样本

对比样本的设计：

**high\_likes\_low\_score**（高赞低分）：

- 选择条件：点赞 > 1000 且 情感分数 < 0.3
- 意义：用户认可（高赞）但算法判负面（低分）
- 用途：发现情感分析误判

**low\_likes\_but\_long**（低赞长文）：

- 选择条件：点赞 < 100 且 字数 > 100
- 意义：内容丰富（长文）但没被发现（低赞）
- 用途：发现被埋没的深度内容

**价值：**

对比样本的价值：

1. 发现算法误判
  - 高赞低分 → 可能是感伤式金句
  - AI 阅读后可以修正判断
2. 发现被埋没的内容
  - 低赞长文 → 可能是深度评论
  - 没有热度≠没有价值
3. 让 AI 有"质疑"能力
  - 不只是接受量化数据
  - 可以挑战算法的判断

## 8.4 为什么是100条阈值

设计过程：

考虑因素：

1. 统计学要求
    - 中心极限定理：样本  $\geq 30$
    - NLP 分析：建议  $\geq 100$
    - 选择 100 作为最低要求
  2. 六维度需求
    - sentiment：需要足够样本计算分布
    - temporal：需要跨年份分布
    - social：需要足够样本计算基尼系数
    - 100条可以满足基础需求
  3. 用户体验
    - 太低：分析质量差
    - 太高：阻挡太多用户
    - 100 是平衡点
-

## 九、典型使用流程

### 完整流程演示

用户: "帮我分析《晴天》的评论区"

Step 1: 搜索 → 确认 → 入库

---

AI 调用: `search_songs_tool("晴天")`

返回: 候选列表

- 1. 晴天 - 周杰伦 | 专辑:叶惠美 | 2003
- 2. 晴天 - 其他歌手...

用户选择: 1

AI 调用: `confirm_song_selection_tool(session_id, choice=1)`

返回: `song_id = "185811"`

AI 调用: `add_song_to_database("185811")`

返回: 入库成功, 热评15条已保存

Step 2: Layer 0 查看数据边界

---

AI 调用: `get_analysis_overview_tool("185811")`

返回:

```
{
  "status": "success",
  "data_boundary": {
    "db_count": 70,
    "api_total": 624567,
    "coverage_rate": 0.0001
  },
  "quality_assessment": {
    "sufficient": false,
    "recommendation": "数据不足, 建议采样"
  }
}
```

AI 决策: 数据不足, 需要采样

Step 3: 采样

---

AI 调用: sample\_comments\_tool("185811", level="standard")

返回:

```
{
  "status": "success",
  "actual": 598,
  "samples": {"hot": 15, "recent": 175, "yearly": 408},
  "coverage": {"years_span": 21, "years_sampled": 10}
}
```

Step 4: 再次 Layer 0 确认

---

AI 调用: get\_analysis\_overview\_tool("185811")

返回:

```
{
  "status": "success",
  "data_boundary": {
    "db_count": 598,
    "api_total": 624567,
    "coverage_rate": 0.001
  },
  "quality_assessment": {
    "sufficient": true,
    "recommendation": "数据充足, 可以继续"
  }
}
```

Step 5: Layer 1 六维度信号

---

AI 调用: get\_analysis\_signals\_tool("185811")

返回:

```
{
  "dimensions": {
    "sentiment": {
      "distribution": {"positive": 0.67, "negative": 0.15, "neutral": 0.18},
      "confidence": 0.6
    },
    "content": {
```



```
    "keywords": ["青春", "回忆", "那年", "单曲循环", "眼泪"]
  },
  "social": {
    "gini_coefficient": 0.78,
    "top10_concentration": 0.43,
    "interpretation": "高度集中, 精英控场型"
  },
  ...
},
"signals": [
  {
    "type": "high_likes_low_sentiment",
    "description": "发现18条高赞但情感分数低的评论",
    "suggestion": "建议查看对比样本验证"
  }
]
}
```

AI 发现: 有18条高赞低分评论, 可能是算法误判

#### Step 6: Layer 2 验证样本

---

AI 调用: `get_analysis_samples_tool("185811")`

返回:

```
{
  "samples": {
    "anchors": {
      "most_liked": [
        {"content": "那年我们十七岁...", "likes": 125000}
      ]
    },
    "contrast": {
      "high_likes_low_score": [
        {
          "content": "我失去了最好的朋友, 这首歌是我们的回忆",
          "likes": 8500,
          "sentiment_score": 0.25,
          "note": "算法可能误判"
        }
      ]
    }
  }
}
```

```
    ]
  }
}
}
```

AI 阅读样本后判断:

- 高赞低分评论都是"感伤式金句"
- 不是负面，是深情怀念
- 实际正面情感比67%更高

## Step 7: 生成报告

---

AI 输出:

### ## 《晴天》评论区分析报告

#### ### 数据说明

- 分析基于 598 条采样评论 (standard级别)
- API 显示总评论 62.4万, 覆盖率 0.1%
- 时间跨度: 2003-2024, 覆盖 10 个年份

#### ### 核心发现

##### 1. \*\*情感特征\*\*

- 算法统计: 正面67%、负面15%、中性18%
- 实际修正: 阅读高赞低分样本后发现, 多数"负面"评论实为感伤式金句
- 真实情感: 正面比例可能高于80%

##### 2. \*\*主题特征\*\*

- 核心关键词: 青春、回忆、那年、单曲循环
- 评论区定位: 青春回忆的情感树洞

##### 3. \*\*社交特征\*\*

- 基尼系数 0.78, TOP10评论占43%点赞
- 特点: 精英控场型, 少数金句获得大量认可

#### ### 代表性评论

> "那年我们十七岁, 现在我们已经三十岁了" — 125000赞

- ### 局限说明
- 采样覆盖率0.1%，可能遗漏部分内容
  - 早期（2003-2010）评论样本较少

# 十、项目价值与反思

## 核心价值

让 AI 在有限样本上做可靠推断，并清楚自己的局限。

展开来说：

- 价值1：解决了"大数据给AI"的问题
- 200万条评论 → 1000条代表性样本
  - 快速（1分钟）且有效

- 价值2：白盒设计建立信任
- AI 知道数据从哪来
  - AI 知道算法的局限
  - AI 可以在报告中诚实说明

- 价值3：AI 可以"质疑"
- 不只是接受量化结果
  - 可以阅读原文验证
  - 可以发现算法盲区

## 学到什么

层面	收获
思维方式	第一性原理：从"用户要什么"出发，不是"技术能做什么"
设计哲学	白盒设计：透明告知比隐藏问题更重要
架构原则	职责边界：工具提供证据，不下结论

层面	收获
产品思维	强制流程：防止 AI 跳过关键步骤
工程实践	API逆向、分层架构、采样策略

## 不足与未来

当前不足：

1. 采样仍有偏差

- 0.1%覆盖率不能代表全部

- 历史评论可能被删除

2. 算法局限明显

- 情感分析误判率高

- 依赖 AI 修正

3. 单歌曲分析

- 还不支持歌曲对比

- 不支持歌单分析

未来方向：

1. 歌曲对比分析

2. 歌单整体画像

3. 更多平台接入（QQ音乐、Spotify）

4. 更好的情感分析模型

---

## 十一、Q&A

**Q1: 采样1000条能代表200万条吗？**

A: 不能完全代表，所以我们：

- 透明告知覆盖率

- 分层采样保证时间覆盖
- 让 AI 在报告中说明局限
- 关键是"知道自己不知道什么"

## **Q2: 为什么工具不直接输出结论?**

A: 因为算法有系统性偏差:

- SnowNLP 对感伤式金句误判
- 工具下结论会误导 AI
- 提供数据+样本, 让 AI 自己判断更可靠

## **Q3: 强制流程检查会不会限制 AI 的灵活性?**

A: 是的, 这是故意的:

- 早期测试发现 AI 会跳过关键步骤
- 强制检查保证分析质量
- 灵活性让位于可靠性

## **Q4: 为什么选择600作为标准采样量?**

A: 经验值, 基于:

- 能覆盖6-10个年份
- 各年份有50-100条样本
- 1分钟内完成
- 六维度分析可靠

## **Q5: cursor参数是怎么发现的?**

A: API逆向分析:

- 抓包分析网易云Web端请求
- 发现weapi接口的cursor是时间戳
- 测试验证可以跳转任意历史时间

## **Q6: 为什么是六个维度?**

A: 基于人对于评论区各个方面的量理解, 后续可以补充优化:

- 情感：听众感受
- 内容：聊什么话题
- 时间：氛围变化
- 结构：评论形态
- 社交：互动模式
- 语言：评论类型

**Q7: 置信度是怎么确定的?**

A: 基于:

- 算法原理 (如SnowNLP训练数据)
- 实际测试误差率
- 同类研究的经验值
- 保守估计, 宁低勿高