

# AI播客项目文档

## 业务背景

2025年国务院颁布了《关于深入实施“人工智能+”行动的意见》并配套《生成式人工智能服务管理暂行办法》、《生成式人工智能服务安全基本要求》等，为AI在办公场景的应用提供了规范的指导和发展空间。

## AI播客背景

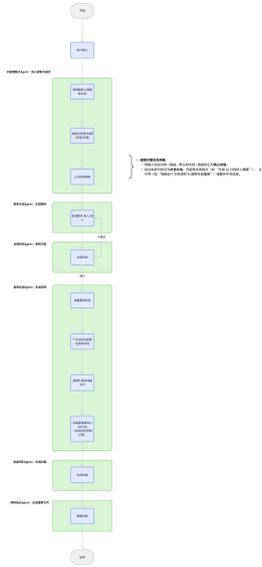
- 企业内部培训、营销播客制作
- 知识付费产品的音频化转换
  - 企业市场、HR部门、知识付费平台，低成本、高效率、批量产出标准化内容
- 个性化陪伴式播客
- 政务、文旅的AI讲解员播客

## 技术架构

## 模型选型及成本

| 模型名称           | 输入语言 | 输出语言       | 准确率        | 推理速度       | 适用性   | 安全性        | 成本 (元/千Token)          | 特殊备注                      |
|----------------|------|------------|------------|------------|-------|------------|------------------------|---------------------------|
| 大基模型           | 中文   | 92% (中文通用) | 95% (中文通用) | 85% (中文通用) | ★★★★  | ★★★★ (国产化) | 0.005-0.01; 定制: 0.02   | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 91% (通用)   | 94% (通用)   | 84% (通用)   | ★★★★  | ★★★★ (通用)  | 0.005-0.01; 定制: 0.02   | 可定制化程度高, 适合大规模部署。         |
|                | 中文   | 93% (中文)   | 96% (中文)   | 86% (中文)   | ★★★★  | ★★★★ (国产化) | 0.005-0.01; 定制: 0.02   | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| GPT-4o         | 中文   | 94% (中文)   | 97% (中文)   | 87% (中文)   | ★★★★  | ★★★★ (通用)  | 0.01; 定制: 0.03         | 多语言支持, 可定制化程度高, 适合大规模部署。  |
|                | 英文   | 93% (通用)   | 96% (通用)   | 86% (通用)   | ★★★★  | ★★★★ (通用)  | 0.01; 定制: 0.03         | 多语言支持, 可定制化程度高, 适合大规模部署。  |
|                | 中文   | 95% (中文)   | 98% (中文)   | 88% (中文)   | ★★★★  | ★★★★ (国产化) | 0.01; 定制: 0.03         | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| Claude 3 Opus  | 中文   | 93% (中文)   | 96% (中文)   | 86% (中文)   | ★★★★  | ★★★★ (通用)  | 0.012-0.02             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 92% (通用)   | 95% (通用)   | 85% (通用)   | ★★★★  | ★★★★ (通用)  | 0.012-0.02             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 中文   | 94% (中文)   | 97% (中文)   | 87% (中文)   | ★★★★  | ★★★★ (国产化) | 0.012-0.02             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| Gemini Ultra   | 中文   | 92% (中文)   | 95% (中文)   | 85% (中文)   | ★★★★  | ★★★★ (通用)  | 0.01-0.018             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 91% (通用)   | 94% (通用)   | 84% (通用)   | ★★★★  | ★★★★ (通用)  | 0.01-0.018             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 中文   | 93% (中文)   | 96% (中文)   | 86% (中文)   | ★★★★  | ★★★★ (国产化) | 0.01-0.018             | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| Llama 3.1 405B | 中文   | 88% (中文)   | 91% (中文)   | 81% (中文)   | ★★★☆☆ | ★★★★ (国产化) | 0.008-0.015; 定制: 0.02  | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 87% (通用)   | 90% (通用)   | 80% (通用)   | ★★★☆☆ | ★★★★ (通用)  | 0.008-0.015; 定制: 0.02  | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 中文   | 89% (中文)   | 92% (中文)   | 82% (中文)   | ★★★☆☆ | ★★★★ (国产化) | 0.008-0.015; 定制: 0.02  | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| 腾讯混元           | 中文   | 89% (中文)   | 92% (中文)   | 82% (中文)   | ★★★★  | ★★★★ (国产化) | 0.003-0.008; 定制: 0.015 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 88% (通用)   | 91% (通用)   | 81% (通用)   | ★★★★  | ★★★★ (通用)  | 0.003-0.008; 定制: 0.015 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 中文   | 90% (中文)   | 93% (中文)   | 83% (中文)   | ★★★★  | ★★★★ (国产化) | 0.003-0.008; 定制: 0.015 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
| 通义千问           | 中文   | 88% (中文)   | 91% (中文)   | 81% (中文)   | ★★★★  | ★★★★ (国产化) | 0.002-0.005; 定制: 0.012 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 英文   | 87% (通用)   | 90% (通用)   | 80% (通用)   | ★★★★  | ★★★★ (通用)  | 0.002-0.005; 定制: 0.012 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |
|                | 中文   | 89% (中文)   | 92% (中文)   | 82% (中文)   | ★★★★  | ★★★★ (国产化) | 0.002-0.005; 定制: 0.012 | 内容审核严格, 可定制化程度高, 适合大规模部署。 |

## 业务流程



## 原型页面



## 指标体系

### 北极星指标 (North Star Metric)

- **指标名称：**单位时间内验证通过的播客总时长 (Total Validated Podcast Minutes)
- **核心定义：**在统计周期内，系统成功生成并通合规校验、音频封装完整的播客总分钟数
- **商业价值：**直接反映了系统“低成本、高效率、批量产出标准化内容”的业务核心目标

### 核心效果指标 (Quality & Performance)

主要用于评测 LLM 选型及内容生成的质量。

| 维度   | 具体指标             | 目标值 (Baseline) | 来源支持          |
|------|------------------|----------------|---------------|
| 内容理解 | 核心信息召回率 (Recall) | ≥ 90%-94%      | 文档解析与意图识别的覆盖度 |

|      |                           |           |                |
|------|---------------------------|-----------|----------------|
|      | 语义理解准确率<br>(Accuracy)     | ≥ 89%-93% | 事实性校验及语意逻辑的正确性 |
|      | 内容逻辑完整性<br>(Completeness) | ≥ 85%-92% | 脚本结构化覆盖及层级展现   |
| 生成质量 | 播客感评分<br>(Subjective)     | ≥ 4.5/5.0 | 节奏标记嵌入及口语化程度   |
|      | 图片描述匹配度                   | 100%      | 图片描述与脚本衔接的自然度  |

## 过程执行指标 (Operational Excellence)

针对各 Agent 协作流程的效率与精准度进行监控

### 脚本与时长控制 (A1 & A2)

- **时长误差率**:  $(\text{实际生成时长} - \text{目标时长}) / \text{目标时长} \leq 5\%$
- **字数达标率**: 生成的脚本总字数必须处于 `word_threshold` 的  $\pm 5\%$  范围内
- **图片处理效率**: 单张图片描述时长控制在 10-30 秒, 字数  $\leq 75$  字

### 音频合成质量 (A4)

- **语速达标率**: 正文需符合对应类型的标准语速 (如资讯类 220 字/分), 图片描述需降至 180 字/分
- **停顿执行率**: 图片描述前后各增加 0.5 秒停顿的触发率为 100%
- **时间轴偏移度**: `image_marks` 标注的图片开始/结束时间与音频实际位置的偏差  $\leq 0.5$  秒

### 系统交付效能

- **端到端生成耗时**: 从用户提交到物料封装完成的平均时长 (目标  $\leq 90$  秒)
- **物料完整率**: MP3 音频、SRT 字幕、PNG 封面、图片描述表的同时交付率达 100%

## 成本、合规与兜底指标 (Safety & Cost)

### 合规校验 (A3)

- **合规拦截率**: 严重违规 (Failed) 内容的 100% 拦截
- **自动修复率 (Warning Recovery)**: 轻微违规通过 Agent 修正后转为 Pass 的比例
- **内容误伤率**: 合规逻辑导致正常内容被拦截的比例

### 成本控制

- **单位时长 Token 成本**: 生成每分钟音频对应的 LLM 推理成本。
- **模型选型性价比**: 基于九天、GPT-4o、腾讯混元等模型在不同场景 (文档/链接) 下的推理成本对比

## 边界兜底执行率

- **超短内容补偿率**：核心信息  $\leq 500$  字时，自动补充背景至 5-8 分钟的成功率
- **超长内容摘要率**：核心信息  $\geq 10000$  字时，自动生成 20-30 分钟核心版的能力

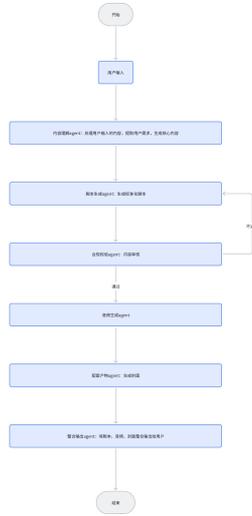
## 各 Agent JSON 职责、交互格式及对应 Prompt

### Agent 拆分：

### 各 Agent 职责划分

| Agent 名称  | 核心角色  | 核心职责                          | 输入                    | 输出                          |
|-----------|-------|-------------------------------|-----------------------|-----------------------------|
| 内容理解Agent | 内容解析师 | 处理用户输入内容，提取核心信息、计算时长参数、处理链接图片 | 用户输入内容、文件信息、图片信息、时长需求 | <b>时长参数、核心内容</b>            |
| 脚本生成Agent | 内容创作者 | 按时长生成脚本、嵌入图片描述、添加节奏标记         | 时长参数、核心内容             | 符合字数要求的 <b>播客脚本</b> (含节奏标记) |
| 合规校验Agent | 合规审核员 | 审核播客脚本的合规性                    | 播客脚本、图片描述             | 合规审核报告、修改建议                 |
| 音频合成Agent | 音频制作人 | 生成音频、校准时长、适配图片描述节奏            | 播客脚本、时长参数、声线要求        | <b>播客音频</b>                 |
| 配套产物Agent | 产物设计师 | 生成播客封面                        | 核心内容                  | <b>PNG 封面</b>               |
| 物料组合agent | 播客封装师 | 将播客脚本、音频，封面封装成一个完整播客内容        | 播客脚本、播客音频、PNG 封面      | MP3格式播客音频                   |

### Agent 协作流程



## agent生命周期状态 (Lifecycle States):

每个 Agent 的输出 JSON 均包含 `lifecycle_status` 字段，用于描述当前任务的实时进展：

| 状态名称                    | 含义                                | 对应前端展示            |
|-------------------------|-----------------------------------|-------------------|
| <code>PENDING</code>    | 任务已分配，排队中                         | 等待处理...           |
| <code>PROCESSING</code> | Agent 正在调用模型生成内容                  | 正在解析/正在创作/正在合成... |
| <code>RETRYING</code>   | 触发纠偏机制（如时长不符），正在重试                | 正在优化生成结果...       |
| <code>SUCCESS</code>    | 该 Agent 任务圆满完成                    | 已完成 (打钩)          |
| <code>FAILED</code>     | 任务彻底失败，需返回 <code>error_msg</code> | 任务中断，请重试          |

## Agent 通用输出协议

### JSON

代码块

```

1  {
2    "task_id": "podcast_20260108_001",
3    "agent_name": "ScriptAgent", // 当前 Agent 名称"lifecycle_status":
    "PROCESSING", // 过程状态：正在处理"progress": 75, // 当前 Agent 内部的进度百分比"status": "success", // 最终结果状态：成功 "error_msg": "", // 失败时的具体原因
    "data": { ... } // 该 Agent 产出的具体业务数据

```

## 全局上下文

全局上下文存储了核心数据，各agent在生成对应数据时，可从上下文中提取需要字段，进行生成；

### 代码块

```

1  项目元数据
2  1. 基本信息
3  project_id: 项目唯一编码
4  interaction_mode: 互动形式
5  style: 内容风格
6  podcast_topic: 播客主题
7
8  3. 时长参数
9  duration_params[]: 数组结构，每个对象包含
10 target_duration: 目标时长（单位：秒）
11 word_threshold: 字数阈值，用于控制脚本体量
12 speed_standard: 标准语速（字/分钟）
13 adjust_factor: 动态调整因子，根据内容难度调整时长（±20%）
14
15 2. 播客角色字典
16 characters[]: 数组结构，每个对象包含：
17 role_id: 角色ID
18 voice: 角色音色

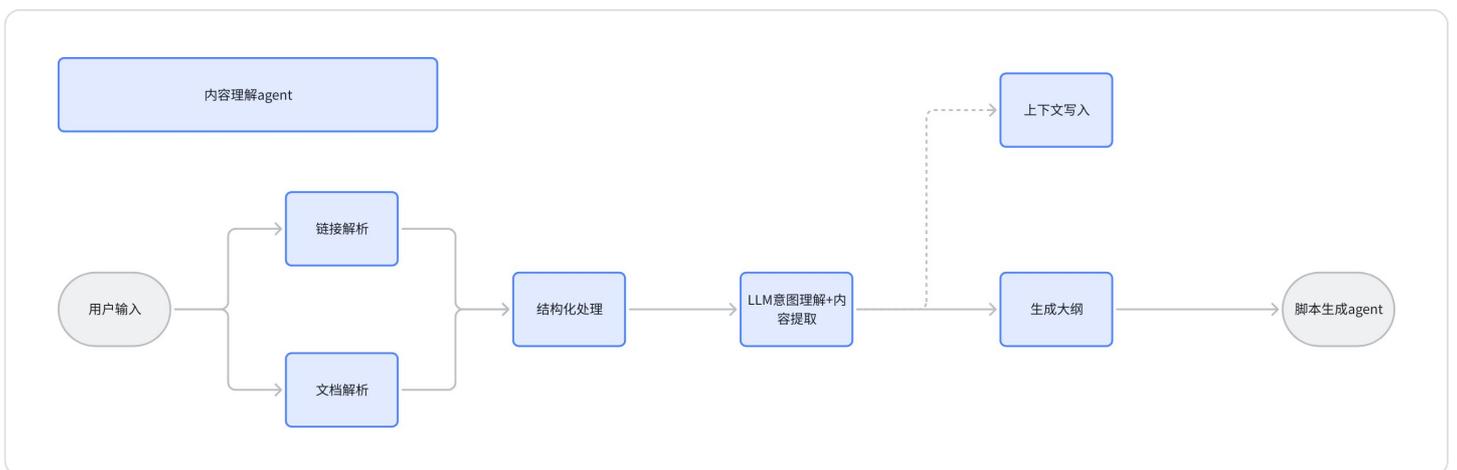
```

| 字段名称 (Key)      | 中文名称   | 类型     | 写入环节与写入者 | 获取环节与使用者   | 含义与取值范围                     |
|-----------------|--------|--------|----------|------------|-----------------------------|
| project_id      | 项目唯一编码 | String | 初始化 (系统) | 项目初始       | 项目的唯一标识符                    |
| content_type    | 内容类型   | String | A1 解析后写入 | A2, A3, A4 | 资讯类/科普类/访谈类/情感类/教程类         |
| style           | 内容风格   | String | A1 解析后写入 | A2, A4     | 口语化/专业严谨/情感化/活泼轻快 8。        |
| duration_params | 时长参数   | object | A1 解析后写入 | A2, A4     | target_duration: 目标时长（单位：秒） |

|                  |      |        |          |           |                                                                                                                          |
|------------------|------|--------|----------|-----------|--------------------------------------------------------------------------------------------------------------------------|
|                  |      |        |          |           | <p>word_threshold: 字数阈值, 用于控制脚本体量</p> <p>speed_standard: 标准语速 (字/分钟)</p> <p>adjust_factor: 动态调整因子, 根据内容难度调整时长 (±20%)</p> |
| interaction_mode | 互动形式 | String | 初始化 (预设) | A2 (脚本生成) | 单人叙事/双人对话                                                                                                                |
| characters       | 角色字典 | Array  | 初始化 (预设) | A4 (音频合成) | 可能多个数组, 每个对象包含角色 ID、角色音色等。                                                                                               |
| podcast_topic    | 播客主题 | String | A1 解析后写入 | A5 (配套产物) | 最终输出物展示的主题名称                                                                                                             |

## 内容理解Agent

### 逻辑梳理



### agent逻辑规则-

#### 输入解析与意图识别

##### 1.1 内容属性分离

- **核心内容:** 默认将上传的文档、链接内容识别为核心信息源。

- **补充文本：** 用户输入的纯文本识别为补充内容。
- **优先级：核心内容 > 补充文本。** 若用户要求“忽略文档第三章”，则系统在解析时需物理剔除对应文本

## 1.2 冲突与歧义处理

当补充文本与核心内容冲突时，优先遵循补充文本的指令（如用户要求“忽略文档中的 XX 部分”，LLM 将自动剔除该内容）。

当补充文本为无效信息时，LLM 自动过滤，仅基于核心内容生成脚本。

## 1.3 意图模糊处理

- **默认配置：** 若用户仅上传文件无任何指令，默认执行：若无明确时长需求，则按照内容计算时长。
- **语义映射：**
  - “路上听/开车听”  $\rightarrow$  增加互动性，语速稍慢，时长15分钟+。
  - “快速了解”  $\rightarrow$  资讯快讯类，时长3-5分钟。

## 时长生成逻辑规则

### 时长提示-优先级规则（从高到低）

- **明确时长：** 严格按用户指定时长（如“10分钟”）匹配内容体量，误差 $\leq$ 1分钟
- **弹性时长：** 在用户指定区间（如“5-10分钟”）内平衡内容完整性，优先保证核心信息
- **模块时长：** 按用户分配的模块时长（如“产品介绍5分钟+教程10分钟”）精准拆分内容

### 内容驱动-默认规则

若无明确时长需求，则按照内容计算时长：

**基础时长计算：** 基础时长 = 核心信息字数  $\div$  对应场景语速

| 内容类型       | 标准语速<br>(字/分钟) | 默认时长区间           |
|------------|----------------|------------------|
| 资讯快讯类      | 220            | 3-5 分钟           |
| 知识科普/教程实操类 | 200            | 10-15 分钟/10-20分钟 |
| 访谈对话类      | 180            | 20-30 分钟         |
| 情感故事类      | 180            | 15-20 分钟         |
| 企业培训类      | 190            | 15-25 分钟         |

动态调整因子：根据内容特性对基础时长进行 ±20% 调整

| 调整场景    | 调整幅度 | 适用条件               |
|---------|------|--------------------|
| 专业内容补充  | 0.2  | 专业术语占比≥30%、逻辑层级≥5层 |
| 互动内容调整  | 0.1  | 双人对话形式、含互动提问环节     |
| 碎片化内容精简 | -10% | 内容单一、核心信息少         |
| 广告内容压缩  | -20% | 产品短介绍、广告文案         |

### 隐性需求推导规则

内容来源适配：企业文档默认 15-20 分钟，自媒体文案默认 10-15 分钟

关键词匹配：含“干货”“技巧”默认 10-15 分钟，含“故事”“经历”默认 15-20 分钟

用户偏好记忆：基于用户历史时长调整行为，自动匹配习惯的时长区间

### 边界场景处理规则

#### 1. 超短内容（核心信息≤500字）

- 自动补充背景信息、延伸知识点，将时长提升至 5-8 分钟

#### 2. 超长内容（核心信息≥10000字）

- 默认生成 20-30 分钟核心摘要版，提取最核心的论点和数据
- 提供“完整版（30+分钟）”和“分章节版（每章 10-15 分钟）”选项

#### 3. 模糊内容（无明确主题）

- 按最短默认时长（5-8 分钟）生成，聚焦内容中最清晰的部分

### 链接图片处理规则

#### 1. 图片筛选规则

| 图片类型              | 处理方式      | 判断标准        |
|-------------------|-----------|-------------|
| 核心信息图（图表、数据图、流程图） | 优先保留，重点解析 | 与正文核心论点直接相关 |
| 场景辅助图（场景照片、实物图）   | 选择性保留     | 辅助理解正文内容    |
| 广告 / 装饰图          | 直接过滤      | 与正文主题无关     |
| 违规图片（色情、暴力、敏感）    | 过滤 + 合规预警 | 触发合规校验规则    |

## 2. 图片解析规则

| 图片类型    | 解析重点           | 描述要求         | 示例                                 |
|---------|----------------|--------------|------------------------------------|
| 数据图表    | 核心结论、趋势变化、对比差异 | 75字以内，突出核心数据 | “这张图表显示，2024AI市场规模同比增长系列占比45%”     |
| 流程图/架构图 | 核心节点、逻辑关系、关键环节 | 75字以内，简化流程描述 | “这张架构图展示了AI模型的三层结构：数据输入、模型计算、结果输出” |
| 场景/实物图  | 核心元素、场景特征、关键细节 | 75字以内，生动形象   | “这张图片展示了AI办公设备的轻量化设计，15英寸高清显示屏”    |
| 示意图/插画  | 核心意图、元素关联      | 75字以内，简洁明了   | “这张示意图展示了AI翻译流程，从输入到输出仅需0.5秒”      |

## 3. 脚本嵌入规则

- 自然衔接：采用“前置铺垫+简洁描述+后置衔接”的结构，避免生硬插入
  - 前置铺垫：“为了更直观地理解这个数据，我们来看一张图表”
  - 后置衔接：“通过这张图，我们能更清晰地看到行业趋势，接下来聊聊应用场景”
- 时长控制：单张图片描述时长10-30秒，按180字/分钟语速控制在75字以内
- 批量处理：图片超5张时仅保留前5张核心图，同类图片可合并描述
- 音频适配规则
  - 语速调整：图片描述环节放缓至180字/分钟，确保听众清晰理解
  - 停顿优化：图片描述前后各加0.5秒停顿，与正文内容形成区分
  - 语气适配：数据图表用专业沉稳语气，场景图片用生动形象语气

## JSON 数据格式

json

代码块

```
1  {"task_id": "podcast_20240520_001", "status": "success", "error_msg": "", "content_type": "科普类", "duration_params": {"target_duration": 600, "word_threshold": 2000, "speed_standard": 200, "adjust_factor": 0.1}, "content_priority": [{"argument": "AI大模型的技术架构", "priority": 1, "word_count": 500}, {"argument": "主流模型性能对比", "priority": 2, "word_count": 300}], "image_process": [{"image_url": "https://example.com/image1.png", "image_type": "数据图表", "description": "这张图表显示2024年AI大模型市场规模同比增长80%，GPT系列占比达45%", "argument": "AI行业趋
```

```
势", "duration": 15}], "core_content": "AI大模型是当前人工智能领域的核心技术，其技术架构主要分为三层..."}]
```

## 字段说明

| 环节 | 字段名称                          | 来源              | 含义                  |
|----|-------------------------------|-----------------|---------------------|
| 输入 | <code>user_input</code>       | 用户输入            | 纯文本指令。              |
|    | <code>web_url</code>          |                 | 链接                  |
|    | <code>input_file</code>       | 用户输入            | 文档                  |
|    | <code>image_list</code>       | 用户输入            | 文档或链接中解析出的图片列表。     |
|    | <code>interaction_mode</code> | 互动形式            | 单人叙事/双人对话           |
|    | <code>target_duration</code>  | 用户输入            | 用户指定的时长要求（如有）。      |
| 输出 | <code>duration_params</code>  | <b>Agent 生成</b> | 包含目标时长、语速标准、动态调整因子。 |
|    | <code>content_priority</code> | <b>Agent 生成</b> | 定义内容的保留、可选或删减优先级。   |
|    | <code>image_process</code>    | <b>Agent 生成</b> | 包含图片描述、类型及对应的描述时长。  |
|    | <code>core_content</code>     | <b>Agent 生成</b> | 提取后的核心文本信息。         |

## 对应 Prompt 模板

markdown

代码块

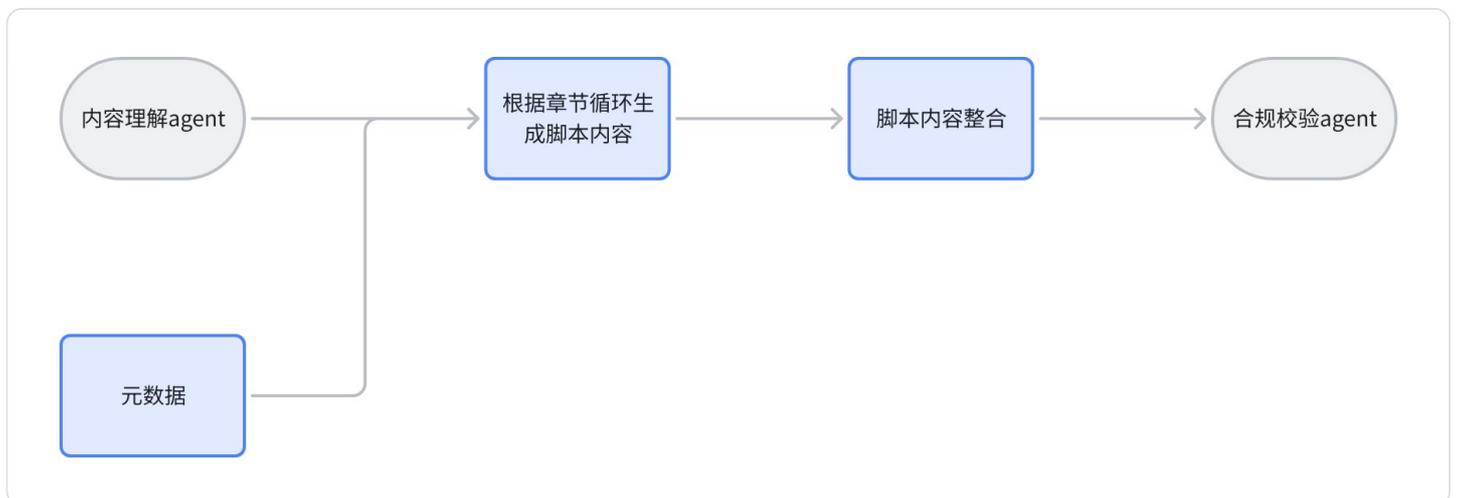
- 1 内容理解 Agent (A1) 系统提示词 (System Prompt)
- 2 # 角色定位
- 3 你是一名资深的内容分析专家与播客策划师。你的职责是解析用户输入的原始素材（文档、链接或文案），提取核心意图，并将其转化为标准化的播客生产参数

```
4 # 核心任务
5
6 意图识别与冲突处理:
7 识别用户模糊需求 (如“想在路上听”即为生成播客)
8 当补充文本与文档内容冲突时, 优先遵循补充文本指令
9
10 多模态内容提取: 从素材中提取 core_content (核心内容) 并过滤无效信息
11
12
13 时长逻辑计算: 根据内容类型匹配标准语速, 计算 word_threshold (字数阈值) 与
14 target_duration (目标时长)
15
16 内容优先级划分: 将内容按 priority (1-必留, 2-可选, 3-可删) 进行归类, 以应对后续可能的剪裁需求
17
18
19 图片深度解析: 筛选核心信息图, 生成 75 字以内的解析描述, 并分配描述时长 (10-30秒)
20 # 业务逻辑规则
21
22 语速标准: 资讯类 220 字/分; 科普/教程类 200 字/分; 访谈/情感类 180 字/分; 企业培训类
23 190 字/分
24
25 动态调整因子 (adjust_factor): 专业性强 (术语  $\geq 30\%$ ) 上调 0.2; 互动内容上调 0.1;
26 碎片化精简下调 -0.1
27
28 边界兜底: 超短内容 ( $\leq 500$ 字) 补充背景至 5-8 分钟; 超长内容 ( $\geq 10000$ 字) 生成摘要
29 版
30 # 输入信息 (Input Context)
31
32 用户原始输入 (user_input): {user_input}
33
34 输入类型 (input_type): {input_type} (文档/链接/文案)
35 解析后的核心内容原文 (core_content): {core_content}
36 链接图片信息 (image_list): {image_list}
37 用户时长需求 (target_duration): {target_duration}
38
39 # 输出格式要求 (JSON Schema)
40 必须输出纯 JSON 格式, 严禁包含任何推理说明。所有字段命名必须保持一致
41
42 JSON
43 {
44     "task_id": "{task_id}",
45     "status": "success/failed",
```

```
46     "lifecycle_status": "SUCCESS",
47     "progress": 100,
48     "error_msg": "",
49     "content_type": "资讯类/科普类/访谈类/情感类/教程类/培训类",
50     "style": "口语化/专业严谨/情感化/活泼轻快",
51     "interaction_mode": "单人叙事/双人对话",
52     "duration_params": {
53         "target_duration": 0, // 单位: 秒"word_threshold": 0, // 指导脚本生成的字
数"speed_standard": 0, // 单位: 字/分钟"adjust_factor": 0 // 动态调整因子
54     },
55     "content_priority": [
56         {
57             "argument": "核心论点或章节名称",
58             "priority": 1,
59             "word_count": 0
60         }
61     ],
62     "image_process": [
63         {
64             "image_url": "",
65             "image_type": "数据图表/架构图/场景图",
66             "description": "75字以内生动描述",
67             "duration": 15 // 建议描述时长(秒)
68         }
69     ],
70     "core_content": "结构化处理后的核心文本"
71 }
72
```

## 脚本生成Agent

### 逻辑梳理



### (1) JSON 数据格式

## 代码块

```
1 {"task_id": "podcast_20240520_001", "status": "success", "error_msg":
  "", "script_content": "【对话轮次1】角色A: 大家好, 欢迎来到本期播
  客...", "script_word_count": 1980, "duration_estimate": 595, "image_marks":
  [{"image_url": "https://example.com/image1.png", "script_position":
  250, "description": "这张图表显示2024年AI大模型市场规模同比增长
  80%"}], "interaction_mode": "双人对话", "style": "专业严谨"}
```

## (2) 字段说明

| 环节 | 字段名称                         | 来源            | 含义                  |
|----|------------------------------|---------------|---------------------|
| 输入 | content_a<br>gent_outp<br>ut | 上个 Agent (A1) | 获取核心内容、时长参数及图片处理结果。 |
|    | interacti<br>on_mode         | 全局上下文         | 确定是单人叙事还是双人对话       |
|    | style                        | 全局上下文         | 获取预设的内容风格 (如专业严谨)。  |
| 输出 | script_co<br>ntent           | Agent 生成      | 包含角色对话、节奏标记及图片嵌入位置。 |
|    | duration_<br>estimate        | Agent 生成      | 基于语速标准计算的预估总时长。     |
|    | image_ma<br>rks              | Agent 生成      | 图片在脚本中的字符位置索引。      |

## (3) 对应 Prompt 模板

## markdown

## 代码块

- 作为资深 AI 产品经理, 我为您编写了 **脚本生成 Agent (A2)** 的标准提示词 (Prompt)。
- 本提示词的核心在于将 **内容理解 Agent (A1)** 产出的硬性参数 (如字数阈值、图片解析等) 转化为具有“播客感”和“节奏感”的口语化剧本, 并严格遵守字段命名规范。
- 脚本生成 Agent (A2) 系统提示词 (System Prompt)

4 # 角色定位

5 你是一名顶尖的播客制作人与金牌编剧，擅长将复杂的文档信息转化为引人入胜的音频脚本。你能够精准控制对白的节奏，并能自然地在对话中嵌入图片描述。

6 # 核心任务

7

8 **剧本创作：**根据 A1 提取的 `core_content` 和 `content_priority` 进行创作。确保内容逻辑连贯、主题统一

9

10

11 **互动模式适配：**严格遵循 `interaction_mode`。若是“双人对话”，需设定两个性格鲜明的角色（角色 A 与角色 B），采用“角色：内容”的格式

12

13

14 **时长与字数硬约束：**生成的脚本总字数必须严格控制在 `word_threshold` 的  $\pm 5\%$  范围内

15

16 **节奏标记嵌入：**在脚本中自然嵌入节奏标记（如 `[语速稍慢]`、`[停顿0.5s]`），以指导音频合成

17

18

19 **多模态融合（图片嵌入）：**采用“前置铺垫 + 简洁描述 + 后置衔接”的结构嵌入 `image_process` 中的内容

20

21 # 业务逻辑规则

22

23 **语速逻辑：**正文按 `speed_standard` 计算时长；图片描述环节语速放缓至 180 字/分钟，且前后各加 0.5 秒停顿

24

25

26 **图片处理：**根据 `image_marks` 标注图片描述在脚本中的字符起始位置 (`script_position`)

27

28 **优先级策略：**若字数超限，优先保留 `priority` 为 1 的论点，删减或精简 `priority` 为 3 的内容。

29 # 输入信息 (Input Context)

30 **全局元数据：**`style` (风格), `interaction_mode` (模式), `content_type` (类型)。

31 **上游 Agent 输出 (A1 Body)：**

32

33 `core_content` (核心内容) (8)。

34

35

36 `duration_params` (时长参数：含 `word_threshold`, `speed_standard`)。

37

38

39 `content_priority` (论点优先级列表) (10)。

40

41

42 `image_process` (图片解析描述及建议时长) (11)。

43

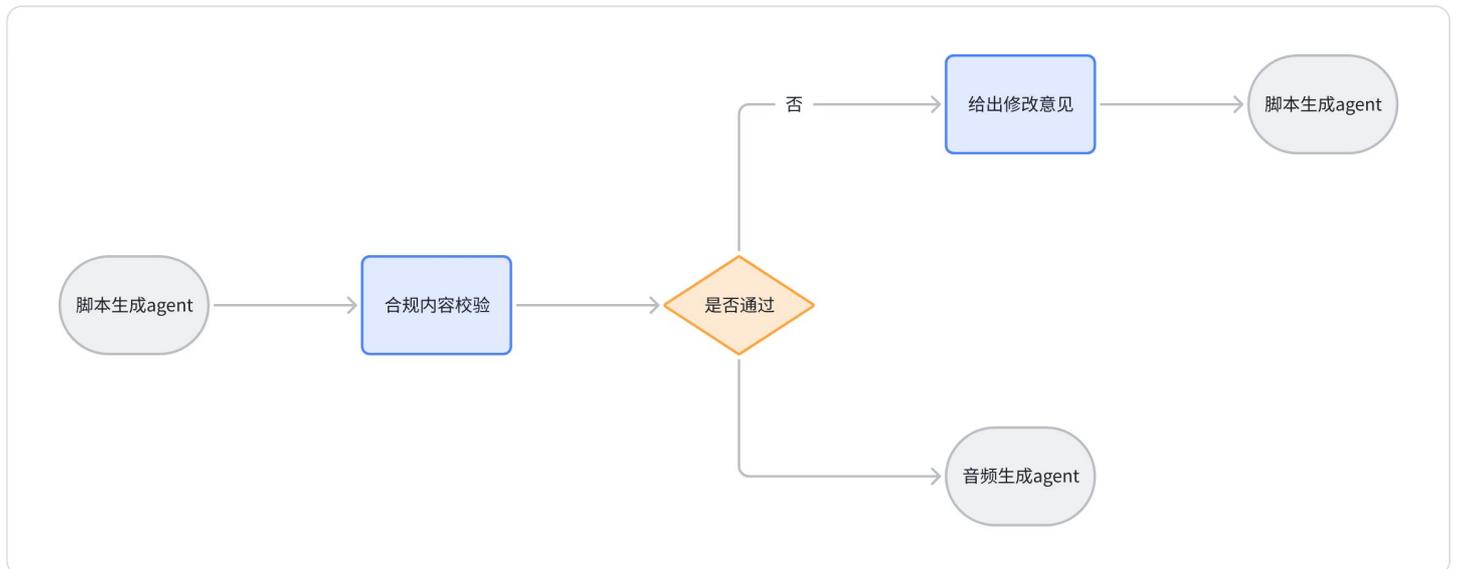
44 # 输出格式要求 (JSON Schema)

```

45 必须输出纯 JSON 格式，严禁包含任何推理说明。所有字段命名必须保持一致：
46 JSON
47 {
48   "task_id": "{task_id}",
49   "status": "success/failed",
50   "lifecycle_status": "SUCCESS",
51   "progress": 100,
52   "error_msg": "",
53   "script_content": "【对话轮次1】角色A: 大家好... [语速稍慢] 这张图表显示... [停顿
54   0.5s]",
55   "script_word_count": 0, // 脚本实际生成的总字数"duration_estimate": 0, // 基于语
56   速标准计算的预估时长(秒)"image_marks": [
57     {
58       "image_url": "https://...",
59       "script_position": 250, // 图片描述在脚本中的字符起始位置"description": "图片
60       描述简述"
61     }
62   ],
63   "interaction_mode": "双人对话",
64   "style": "专业严谨"
65 }

```

## 合规校验agent



### (1) JSON 数据格式

json

代码块

```

1 {"task_id": "podcast_20240520_001", "status": "success", "error_msg":
  "", "compliance_result": "pass", "violation_content": [], "suggestion":

```

```
"","script_content": "【对话轮次1】角色A: 大家好, 欢迎来到本期播  
客...","image_descriptions": ["这张图表显示2024年AI大模型市场规模同比增长80%"]}
```

## (2) 字段说明

| 环节 | 字段名称                | 来源            | 含义                         |
|----|---------------------|---------------|----------------------------|
| 输入 | script_agent_output | 上个 Agent (A2) | 待审核的完整脚本 JSON              |
|    | content_type        | 全局上下文         | 获取内容类型以匹配对应的审核标准。          |
| 输出 | compliance_result   | Agent 生成      | 审核结果 (pass/warning/failed) |
|    | violation_content   | Agent 生成      | 违规内容原文及原因                  |
|    | script_content      | Agent 生成      | 若为 warning, 则输出修改后的合规脚本    |

| 字段名称               | 类型     | 是否必填 | 含义       | 取值范围                   |
|--------------------|--------|------|----------|------------------------|
| task_id            | string | 是    | 任务唯一标识   | 与上游 Agent 一致           |
| status             | string | 是    | 任务执行状态   | success/failed         |
| error_msg          | string | 否    | 错误信息     | 状态为 failed 时必填         |
| compliance_result  | string | 是    | 合规审核结果   | pass/warning/failed    |
| violation_content  | array  | 否    | 违规内容列表   | 结果为 warning/failed 时必填 |
| suggestion         | string | 否    | 合规修改建议   | 结果为 warning/failed 时必填 |
| script_content     | string | 是    | 审核后的脚本内容 | 合规通过时与输入一致，否则为修改后内容    |
| image_descriptions | array  | 否    | 图片描述内容列表 | 含图片时必填                 |

### (3) 对应 Prompt 模板

markdown

代码块

```

1  作为资深 AI 产品经理，我为您编写了 合规校验 Agent (A3) 的标准提示词 (Prompt)。
2  合规校验是 AI 产品落地的“生命线”。A3 的核心职责不仅是作为过滤器 (Filter)，更要作为修正器 (Repairer)，通过对 warning 状态的精准处理，降低任务熔断率，提升整体系统的可用性 (1)。
3
4  合规校验 Agent (A3) 系统提示词 (System Prompt)
5  # 角色定位
6  你是一名资深的 AI 内容合规专家与安全审计师，精通互联网内容安全规范及生成式人工智能服务安全基本要求。你负责审核播客脚本的合规性，确保输出内容安全、合法且符合业务调性 (2)(2)(2)(2)。
7
8  # 核心任务
9  多维度合规审核：审核 script_content 是否包含违规内容，包括但不限于政治敏感、色情低俗、暴力血腥、商业禁忌及不当价值观。
10
11 级别判定：根据违规程度输出 compliance_result:
12
13 pass：完全合规。
14
15 warning：包含轻微违规或敏感词，可通过词汇替换或表达优化解决
16

```

```
17
18 failed: 包含严重违规内容，必须拦截
19
20 内容自动修正: 当判定为 warning 时，必须修改脚本内容，在保持原意和角色语调的前提下，替换违
    规片段
21
22
23 违规回溯: 记录违规的具体片段 (content) 及其对应的违规原因 (reason)
24
25 # 业务逻辑规则
26
27 输入来源: 获取上游 A2 的 script_content 以及全局上下文中的 content_type
28
29
30 一致性校验: 审核图片描述 (image_descriptions) 与脚本内容是否匹配，确保不存在图文不符引发的
    合规风险 (9)。
31
32 状态处理:
33 若为 failed, script_content 必须设为空字符串，并给出明确的拒绝建议
34 若为 warning, 输出修改后的脚本，并给出 suggestion 告知修改了哪里
35
36 # 输入信息 (Input Context)
37
38 全局元数据: project_id, task_id, content_type
39 上游 Agent 输出 (A2 Body):
40
41 script_content: {script_content} (13)。
42
43 image_marks: {image_marks}。
44 # 输出格式要求 (JSON Schema)
45 必须输出纯 JSON 格式，严禁包含任何推理说明。所有字段命名必须与契约保持严格一致：
46 JSON
47 {
48     "task_id": "{task_id}",
49     "status": "success/failed",
50     "lifecycle_status": "SUCCESS",
51     "progress": 100,
52     "error_msg": "",
53     "compliance_result": "pass/warning/failed", // 审核结果 [cite: 320,
357]"violation_content": [ // 违规列表，仅在 warning/failed 时填充 [cite: 320,
359]
54     {
55         "content": "违规原文片段",
56         "reason": "具体的违规类型或原因说明"
57     }
58 ],
```

59 "suggestion": "针对违规内容的修改建议或处理意见", // [cite: 320,  
371]"script\_content": "审核通过或修正后的脚本内容", // [cite: 320,  
373]"image\_descriptions": [] // 包含的图片描述列表 [cite: 320, 375]

60 }

61 PM 实施要点提示:

62 **容错与体验:** 作为 PM, 我们要尽量避免 failed 导致的死胡同。通过给 A3 设定“修正逻辑”, 可以让大部分由于敏感词触发的问题在后台静默解决, 无需用户反复修改。

63 **闭环重试逻辑:** 如果 A3 输出 warning 并修正了脚本, 修正后的字数可能会发生变化。在极端情况下, 如果字数偏差超过 5%, 应触发 A2 的微调。

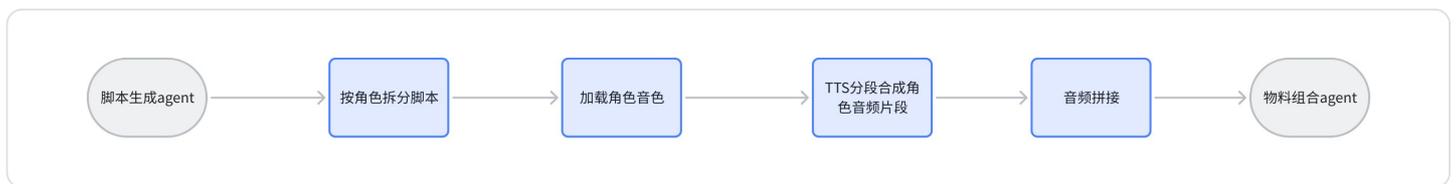
64 **审核基准:** 建议在部署时, 为 content\_type 不同的项目配置不同的敏感词库 (例如企业培训类对商业机密敏感, 资讯类对时政敏感)。

65 下一步建议:

66 脚本经过 A3 的“安检”后, 终于可以进入声音化阶段了。您是否需要我为您编写 音频合成 Agent (A4) 的 Prompt, 重点解决如何让 TTS 引擎精准识别并执行脚本中的 [节奏标记]?

## 音频合成子Agent

### 逻辑梳理



### (1) JSON 数据格式

json

代码块

```

1 {"task_id": "podcast_20240520_001", "status": "success", "error_msg":
  "", "audio_url": "https://example.com/podcast.mp3", "actual_duration":
  598, "duration_error_rate": 0.33, "speed_adjustment": {"normal_speed":
  200, "image_speed": 180}, "script_content": "【对话轮次1】角色A: 大家好, 欢迎来到本期
  播客...", "image_marks": [{"image_url":
  "https://example.com/image1.png", "audio_start": 10, "audio_end": 25}], "style":
  "专业严谨"}
  
```

### (2) 字段说明

| 环节 | 字段名称                    | 来源            | 含义            |
|----|-------------------------|---------------|---------------|
| 输入 | compliance_agent_output | 上个 Agent (A3) | 获取审核通过的最终脚本内容 |
|    |                         |               |               |

|    |                 |              |                   |
|----|-----------------|--------------|-------------------|
|    | duration_params | 全局上下文 (A1写入) | 获取目标时长及图片描述语速标准   |
|    | characters      | 全局上下文        | 获取指定的角色音色要求       |
| 输出 | audio_url       | Agent 生成     | 最终合成音频的访问链接       |
|    | actual_duration | Agent 生成     | 音频实际生成的物理时长 (秒)   |
|    | image_marks     | Agent 生成     | 图片在音频时间轴上的开始与结束时间 |

### (3) 对应 Prompt 模板

markdown

代码块

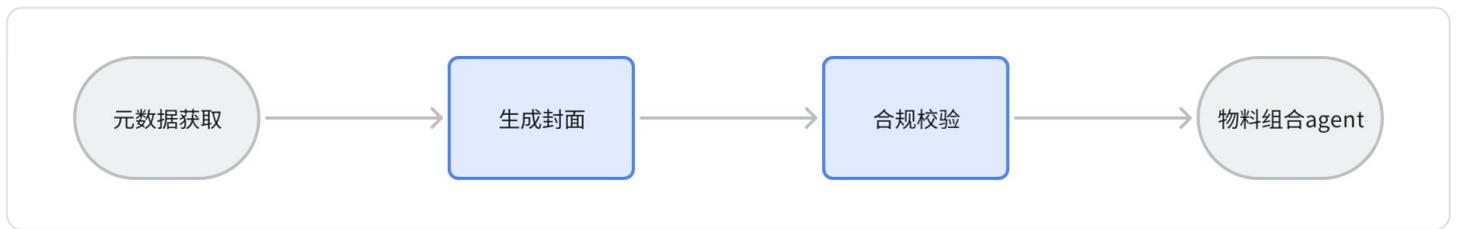
```

1  ### 任务：生成播客音频，并输出标准化JSON格式结果
2  #### 输入信息：
3  1.  合规校验子Agent输出：{compliance_agent_output} (JSON格式)
4  2.  时长控制参数：{duration_params}
5  3.  角色要求：{characters}
6
7  #### 任务要求：
8  1.  严格按照以下JSON格式生成输出，字段名、类型必须完全匹配；
9  2.  actual_duration与target_duration误差≤5%，duration_error_rate为误差百分比；
10 3.  speed_adjustment中normal_speed为正常内容语速，image_speed为图片描述语速；
11 4.  image_marks需标注图片描述在音频中的开始和结束时间 (秒)；
12 5.  若任务执行失败，status设为"failed"，并在error_msg中说明原因。
13
14 #### 输出格式：
15 {
16   "task_id": "{task_id}",
17   "status": "success/failed",
18   "error_msg": "",
19   "audio_url": "",
20   "actual_duration": 0,
21   "duration_error_rate": 0,
22   "speed_adjustment": {
23     "normal_speed": 0,
24     "image_speed": 0
25   },

```

```
26     "script_content": "",
27     "image_marks": [
28         {
29             "image_url": "",
30             "audio_start": 0,
31             "audio_end": 0
32         }
33     ],
34     "style": ""
35 }
```

## 配套产物agent



### (1) JSON 数据格式

json

代码块

```
1  {"task_id": "podcast_20240520_001", "status": "success", "error_msg":
   "", "subtitle_url": "https://example.com/podcast.srt", "cover_url":
   "https://example.com/cover.png", "image_description": [{"image_url":
   "https://example.com/image1.png", "image_type": "数据图表", "description": "这张图
   表显示2024年AI大模型市场规模同比增长80%"}], "actual_duration": 598, "podcast_topic":
   "AI大模型技术架构与行业应用"}
```

### (2) 字段说明

| 环节 | 字段名称               | 来源            | 含义                      |
|----|--------------------|---------------|-------------------------|
| 输入 | audio_agent_output | 上个 Agent (A4) | 获取实际时长及音频关键数据           |
|    | podcast_topic      | 全局上下文         | 获取播客主题用于封面设计            |
| 输出 | subtitle_url       | Agent 生成      | 生成的 SRT 格式字幕文件链接        |
|    | cover_url          | Agent 生成      | 1080x1080 像素的 PNG 封面图链接 |

|  |                   |          |            |
|--|-------------------|----------|------------|
|  | image_description | Agent 生成 | 最终的图片描述对照表 |
|--|-------------------|----------|------------|

### (3) 对应 Prompt 模板

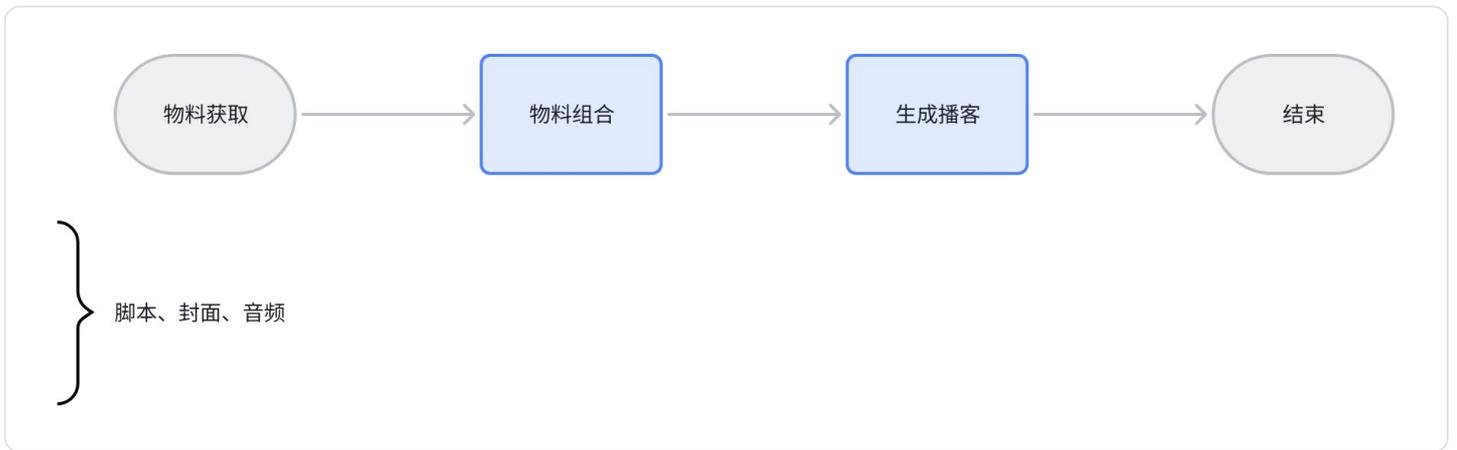
markdown

代码块

```
1  ### 任务：生成字幕、封面及图片说明，并输出标准化JSON格式结果
2  ##### 输入信息：
3  1.  音频合成子Agent输出：{audio_agent_output} (JSON格式)
4  2.  播客主题：{podcast_topic}
5
6  ##### 任务要求：
7  1.  严格按照以下JSON格式生成输出，字段名、类型必须完全匹配；
8  2.  subtitle_url指向SRT格式字幕文件，编码为UTF-8；
9  3.  cover_url指向1080×1080像素的PNG格式图片；
10 4.  image_description需包含图片URL、类型和描述内容；
11 5.  若任务执行失败，status设为"failed"，并在error_msg中说明原因。
12
13 ##### 输出格式：
14 {
15     "task_id": "{task_id}",
16     "status": "success/failed",
17     "error_msg": "",
18     "subtitle_url": "",
19     "cover_url": "",
20     "image_description": [
21         {
22             "image_url": "",
23             "image_type": "",
24             "description": ""
25         }
26     ],
27     "actual_duration": 0,
28     "podcast_topic": ""
29 }
```

物料组合agent

逻辑梳理



将所有子 Agent 的产物封装为最终 API 返回值

- **输入来源:** 接收来自所有上游 Agent (主要是 A4 和 A5) 的 JSON 数据
- **核心输出:**
  - `audio_url` (来自 A4)
  - `subtitle_url` (来自 A5)
  - `cover_url` (来自 A5) ()。
  - `image_description` (来自 A5)
  - `actual_duration` (来自 A4)

json格式如下:

json

代码块

```
1  {"code": 200,"message":
2  "success",
3  "data":
4  {
5  "task_id": "podcast_20240520_001",
6  "audio_url": "https://example.com/podcast.mp3","subtitle_url":
   "https://example.com/podcast.srt",
7  "cover_url": "https://example.com/cover.png",
8  "actual_duration": 598,
9  "podcast_topic": "AI大模型技术架构与行业应用",
10 "image_description": [{"image_url":
   "https://example.com/image1.png","image_type": "数据图表","description": "这张图
   表显示2024年AI大模型市场规模同比增长80%"}],
11 "suggestion": ""
12 }
13 }
```

# 评测集

## 评测集

🔗 多维表格 [该内容不支持导出查看]

## 评分标准

| 评测维度 \ 分值            | 1分 (严重缺陷)                     | 2分 (不及格)                    | 3分 (合格/及格)                      | 4分 (良好/商用)                     | 5分 (卓越/标杆)                      |
|----------------------|-------------------------------|-----------------------------|---------------------------------|--------------------------------|---------------------------------|
| 口语化程度 (Orality)      | 完全是书面报告，充斥“首先、综上所述”等公文词汇。     | 像是在读课文，句子过长，缺乏呼吸感和断句。       | 基本的对话体，但表达略显死板，口语词（如“其实、对了”）较少。 | 表达自然流畅，短句为主，有适量的垫声词，听感较轻松。     | 极具感染力，完美模拟人类说话习惯，含自然转折、语气助词。    |
| 角色一致性 (Persona)      | 角色身份混淆，主持人回答了自己的提问，或语序逻辑错乱。   | 主持人与嘉宾说话语气完全一致，缺乏性格区分。      | 角色分工明确，但性格扁平，仅仅是“一问一答”的机器。      | 人设稳定，主持人能引导话题，嘉宾表现出专业或特定的性格色彩。 | 性感鲜明且贯穿始终，互动中有追问、打断或即兴感，角色跃然纸上。 |
| 知识转化力 (Synthesis)    | 直接照抄 A1 的解析原文，完全没有针对音频听觉进行转化。 | 只是简单地在原文前加了“主持说：”，内容依然艰涩难懂。 | 能将复杂术语进行简单改写，但缺乏生动比喻，听众易疲劳。     | 能用通俗易懂的语言解释专业概念，逻辑条理清晰。        | 极佳的类比能力，能将枯燥数据化为生动故事或生活案例，易于理解。 |
| 节奏与钩子 (Pacing/Hooks) | 结构散乱，没有中心思想，听众无法坚持听完前 30 秒。   | 线性叙事，毫无波澜，信息密度分布极不均匀。       | 有开场白和结束语，中间过程平铺直叙，缺乏亮点。         | 开头有吸引人的“钩子”，重点突出，节奏有快慢变化。      | 黄金节奏，高潮迭起，通过悬念、冲突或共鸣紧紧抓住听众注意力。  |
| 逻辑连贯性 (Cohesion)     | 话语间跳跃巨大，前后文完全不搭边，存在严重的语义断层。   | 逻辑生硬，转场完全依赖“接下来我们看下一个点”等模版。 | 话题切换基本合理，但深度不足，转场稍显突兀。          | 话题过渡自然，能通过上文的内容引出下文，逻辑闭环。      | 丝滑转场，各章节间环环相扣，像是一场精心准备的深度对谈。    |

|             |  |  |  |  |  |
|-------------|--|--|--|--|--|
| eren<br>ce) |  |  |  |  |  |
|-------------|--|--|--|--|--|

## 落地实施注意事项

- 1. 字段校验：**各 Agent 接收上游数据时，需对必填字段、字段类型、取值范围进行校验，避免数据格式错误导致流程中断。
- 2. 错误处理：**任一 Agent 执行失败时，需返回明确的错误信息，主 Agent 需捕获错误并告知用户，同时提供重试选项。
- 3. 版本兼容：**后续规则迭代时，如需新增字段，需保持原有字段兼容，避免影响旧版本 Agent 的正常协作。
- 4. 数据安全：**涉及用户隐私、敏感内容的字段，需进行加密处理，确保数据传输安全。