

AI语音驱动2D数字人直播系统产品方案

1. 产品需求文档（PRD）

产品目标

本项目旨在打造一个利用AI技术驱动的虚拟数字人直播系统，以降低直播门槛、提升内容创作效率和丰富互动形式。通过AI语音合成驱动2D数字人形象进行直播，可实现“无人直播”或辅助真人主播，使企业或个人能够7×24小时持续输出内容，降低人力成本¹²。系统一期聚焦于支持2D数字人的口播直播，涵盖预先录制内容的播放和直播过程中的智能互动；二期将在此基础上迭代支持3D数字人，更加逼真生动的形象表现。最终产品目标包括：

- **降低直播成本**：无需真人主播出镜和专业摄制设备，依托AI数字人即可开展直播¹²。
- **提升内容生产效率**：通过AI自动生成直播脚本、语音及数字人动画，实现快速创建直播内容，减少人工脚本编写和视频录制时间³。
- **增强互动体验**：数字人可与观众实时互动，回答提问或根据弹幕反馈调整内容，让观众获得新奇的互动直播体验。
- **多平台覆盖**：通过OBS推流方式对接抖音、快手、小红书等主流平台，实现一处创作、多平台分发，扩大影响力。

用户角色

本系统的目标用户主要包括： - **内容创作者/运营**：使用本系统制作数字人直播的个人主播、电商商家运营或媒体内容创作者。他们负责设置数字人形象、准备或审核直播内容（如脚本、问答库），并监控直播效果。 - **观众用户**：在抖音、快手、小红书等平台观看数字人直播的用户。他们通过评论/弹幕与数字人互动（提问、下单咨询等），间接参与本系统的互动模块。观众虽非直接使用本系统，但其体验高度依赖系统输出的直播质量。 - **系统管理员（可选）**：维护数字人直播平台运行的技术人员，负责内容审核、系统参数配置、数据管理等，保障系统稳定合规运行。

功能模块

根据以上目标，我们规划以下核心功能模块：

1. **数字人形象生成模块**：提供数字人/avatar的创建与管理功能，支持**自动生成**和**半自动生成**两种模式。自动生成指利用AI根据少量输入（如几张照片或文本描述）快速生成逼真的虚拟形象⁴；半自动生成则提供预设的形象模板，由用户选择基础形象并自定义外观细节（发型、服装、配色等）来打造数字人。每个数字人形象可绑定特定的人设（角色姓名、风格定位）和**声音**（可从AI语音中选择音色或自定义克隆）以保持风格统一⁵⁶。
2. **内容脚本编辑模块**：帮助用户准备直播内容，支持**预录脚本**和**问答库**两类内容管理。预录脚本用于系统在直播时自动播报，包括台本台词、解说稿等；问答库则用于互动直播时的知识储备，包括常见观众问题及对应回答。当用户不提供脚本时，系统还能基于简单输入自动生成直播文案和稿本。例如，用户仅需提供产品链接或卖点描述，AI便可自动撰写完整的直播讲解脚本⁷。

3. **直播场景配置模块**：用于设置直播画面效果，包括场景背景、更换PPT/商品展示画面等（类似“配装修”环节⁸）。一期2D数字人阶段，可提供多套**2D直播场景模板**（如带货场景、才艺展示场景等），用户可选择并替换其中的元素（品牌Logo、产品图片等）。二期3D数字人阶段，将支持更丰富的虚拟场景和道具，营造逼真的直播间环境。
4. **AI语音合成与驱动模块**：本系统的核心模块，负责将文本转化为数字人直播时的语音和相应的动作表情。其子功能包括：
 5. **语音合成（TTS）**：采用深度学习语音合成技术，将脚本文本转为自然流畅的语音⁹。支持多种音色选择，包含普通话男声、女声，以及英文等多语言音色，以满足不同直播人设和语言需求¹⁰。语音合成要求发音清晰自然，语调情感符合语境。
 6. **嘴型/表情动画生成**：根据生成的语音，驱动2D数字人实时做出口型同步的说话动画¹¹。2D数字人采用预先训练的音频驱动模型，将语音特征映射为口型变化、头部动作和适当的表情变化，实现“对口型”的动画效果。可引用先进的唇形同步技术（如Wav2Lip或MuseTalk模型）以提高实时性和准确度¹²。该模块保证数字人的唇动与语音高度同步，延迟控制在可察觉范围之下（<0.2秒）。
 7. **直播互动模块**：用于支撑互动直播场景，在观众通过弹幕/评论提问时，数字人能够智能回应。其功能涵盖：
 8. **评论获取与解析**：集成各平台直播间弹幕/评论获取接口（或通过网页端解析弹幕流），实时收集观众提问内容。采用NLP技术对文本评论进行意图识别和问题解析¹³（如识别询问产品价格、库存等）。
 9. **回答生成**：对于匹配问答库的问题，从预设答案中调取适当回答；对于未收录的问题，可调用大语言模型（LLM）生成回复¹¹。生成的回答再经语音合成与动画模块驱动数字人回答。整套流程确保在几秒内完成，使数字人能够以接近实时的速度与观众对话¹⁴¹¹。
 10. **互动控制**：提供人工介入接口，在半自动模式下运营可以实时编辑或审核AI给出的回复，再由数字人播报，从而在人机协同下保证内容准确性和安全性。
11. **直播发布与推流模块**：实现将数字人直播画面和音频推送到目标平台。由于不采用平台官方SDK，我们通过**OBS推流**方式对接：
 12. 系统与OBS Studio集成，利用OBS的虚拟摄像头功能将数字人直播画面作为摄像头信号输入各平台官方直播工具¹⁵。例如，对于抖音，用户可启动OBS虚拟摄像头，将其设置为抖音“直播伴侣”软件的视频来源，从而绕过抖音对第三方推流的限制¹⁵。
 13. 或者直接使用OBS/FFmpeg推流协议（RTMP）将视频流发送至平台服务器。对于开放推流码的平台（如部分电商平台），系统可让用户配置RTMP地址和串流密钥，然后程序内置推流引擎将数字人音视频流推送出去¹⁶¹⁷。需要注意抖音、视频号等平台默认不公开推流码，强行获取可能导致封禁¹⁸¹⁹，因此推荐采用官方直播软件配合虚拟摄像头的方案确保合规。
 14. 推流模块同时提供**多平台同步直播**支持。通过安装OBS多平台推流插件或启动多路推流进程，可将同一数字人直播同时发送到多个平台账号²⁰（例如抖音和快手），扩大覆盖受众。在多平台直播时注意硬件性能和带宽要求。

15. **数据和管理模块**：提供必要的后台管理和数据记录功能，包括用户账号管理、数字人和脚本内容的存储、直播日志记录等。支持查看直播间关键数据（观看人数、互动次数等）以及内容审核接口，确保生成的内容不违规（可集成敏感词过滤和审核AI）。

交互流程

结合上述模块，典型的用户交互流程如下：

(1) 数字人创建流程：用户首先登录系统，进入“数字人形象”管理。用户选择“新建数字人”，在弹出的界面选择**自动生成或自定义**：

- **自动生成**：上传几张照片或输入文字描述（如“20多岁女性，亲和力风格”），系统调用AI模型生成对应外观的数字人脸部形象和身体立绘⁴。用户确认或让系统再次生成直到满意为止。
- **半自动自定义**：从系统预置的人物形象库中选择一个模板（如某卡通风格女性），然后手动调整肤色、发型、服装等选项，实时预览效果。确定形象后，用户为数字人命名，并选择其语音音色（可试听AI语音样例）。完成后保存数字人档案。

(2) 直播内容配置流程：用户进入“内容脚本”模块，新建一场直播内容。主要步骤包括：

- 选择直播使用的数字人（从已创建列表中选择对应人设）。
- 选择直播模式：**预设脚本模式**或**互动问答模式**，或者两者结合。
- 在预设脚本模式下，用户可以直接撰写直播台本文本，或填写关键点让系统自动生成完整脚本⁷。系统支持将长文案分段编辑，每段关联展示的素材（如产品图、PPT页面）。
- 在互动问答模式下，用户可维护FAQ问答库：输入常见问题及标准回答。也可以启用AI接入，使数字人能基于训练的大模型自动回答未知问题¹¹。
- 选择直播场景模板并配置元素：如上传直播背景图、更换展示素材等，以完善直播画面设计。
- 预览彩排：用户可点击“预览”，系统将模拟播放一小段脚本，由数字人出镜朗读，用户检查语音和动作是否协调，内容是否正确。如有问题可编辑脚本或调整语速、重音等参数，然后再次预览直到满意为止。

(3) 直播执行流程：当内容准备完成后，用户进入“直播发布”环节：

- **填写直播信息**：选择目标平台并登录对应平台账号，设置直播标题、封面等（通过各平台官方工具或网页完成这一步骤）。
- **连接推流**：如果使用OBS虚拟摄像头方案，用户需启动OBS并开启虚拟摄像，再在平台官方直播软件中选择该虚拟摄像头为输入¹⁵。如果使用RTMP推流方案，则在系统中填写获取的平台推流地址和码¹⁶。
- **开始直播**：点击系统中的“开始直播”按钮。此时系统开始按照脚本顺序将文本交给语音合成模块生成语音，并同步驱动数字人动画，将音视频流发送至OBS或RTMP推流端口。数字人出现在目标平台的直播间中开始播报内容。
- **直播过程中**：
 - **如果是预录脚本模式**：数字人将按既定脚本顺序逐段播报内容，期间实时监测弹幕。如果检测到观众提问且已启用互动功能，系统将暂停播报脚本内容，优先回答观众问题，然后再恢复脚本（或跳到下一个相关段落）。
 - **如果是互动直播模式**：数字人主要根据观众提问进行讲解。每当有新问题时，系统经过ASR/NLP（若语音输入）或直接NLP处理文字问题，再通过LLM或问答库生成答案，立刻语音合成并由数字人答复¹¹。运营人员可在后台干预答案或通过提示器给数字人下达指令（如引导话题）。
- **结束直播**：用户手动点击“结束”。系统停止推流，并保存本次直播的日志数据（包括观众提问和AI回答内容，用于改进问答库）。平台直播间关闭后，系统提供直播回放视频的下载或保存。

上述流程确保了用户从创建数字人、准备内容到实际开播的全链路体验。其中一期重点支持预设脚本播放为主，搭配有限的问答互动；二期随着3D数字人和更强AI接入，可实现更丰富的实时对话和表情动作，让直播更加栩栩如生。

验收标准

为确保产品功能达到预期，我们制定以下验收标准：

- **数字人形象生成**：自动生成的数字人形象需在外观上满足用户要求，与输入描述或照片高度相似，面部特征清晰，不出现畸形瑕疵⁴。半自动生成的编辑流程顺畅，可实时预览，无明显界面卡顿。验收标准：随机抽取10次自动生成结果，其中80%以上用户反馈形象质量令人满意；手动定制过程出现重大UI问题为0。
- **语音合成与口型同步**：数字人语音的自然度和清晰度需达到可商业播报水准，语音与对应的口型动作严丝合缝同步¹¹¹²。验收标准：在10段不同内容的测试脚本中，数字人语音自然度平均MOS评分 ≥ 4.0 （满分5）；口型同步延迟 < 0.2 秒且没有明显错位现象¹²。
- **预录内容播报**：数字人能够完整连续地播报至少1小时的脚本内容而不中断，期间音视频流畅无卡顿，表情动作协调无异常定格。验收标准：运行1小时连续脚本测试，音视频无丢帧、无崩溃重启；内容播报正确率100%（无漏读跳读）。
- **互动问答**：在模拟的直播互动场景中，数字人对观众提问的响应正确率和及时性达到要求。验收标准：构造100个FAQ问题集，数字人正确回答率 $\geq 90\%$ （匹配问答库的问题100%准确，开放问题AI回答需逻辑合理无跑题）；对于观众提问，平均响应时间 < 5 秒¹⁴。在半自动模式下，人工干预能够在3秒内插入或修改回答，不影响直播流顺畅。
- **多平台推流**：系统通过OBS成功接入抖音、快手、小红书等平台进行直播测试，各平台直播画面和声音正常。在遵循平台规则前提下，直播至少持续30分钟不断流、不卡顿。验收标准：在上述平台各进行不少于3次试播，推流连接成功率100%，平均上行码率稳定且无异常中断；采用OBS虚拟摄像头接入的平台（如抖音）经验证未触发违规警告¹⁸¹⁹。
- **易用性和稳定性**：产品界面友好，主要功能操作步骤不超过3步即可完成，首次使用有引导提示。系统在常规硬件环境下（如主流笔记本电脑+i7 CPU）运行流畅。验收标准：新手用户不看文档可在15分钟内完成一次完整的直播流程设置并开播；连续运行4小时系统无内存泄漏或崩溃。

通过以上验收指标的逐项测试，确保产品在功能完整性、AI效果和系统稳定性方面达到发布要求。

2. 项目整体规划

阶段划分

项目将按照**两个主要阶段**逐步实施，采用里程碑交付与验收：

- **第一阶段：AI驱动2D数字人直播系统 (预计3个月)**
 - 里程碑1.1 - 技术预研与原型搭建（~0.5个月）：完成详细需求分析和选型论证，包括语音合成、唇动驱动方案、OBS推流方案调研。产出技术方案文档和原型计划。
 - 里程碑1.2 - 核心功能开发完毕（~1.5个月）：实现2D数字人基础功能模块，包括数字人创建、脚本编辑、TTS合成、口型动画驱动和OBS推流集成。内部联调一个端到端流程，验证数字人能朗读脚本并成功推流到测试平台。
 - 里程碑1.3 - 互动功能与UI完善（~0.5个月）：开发弹幕获取和简单问答响应功能，完善前端交互界面（数

字人配置、内容编辑、预览、启动/停止直播等界面）。进行功能测试和UI优化，确保易用性。

里程碑1.4 - 内测验收（~0.5个月）：邀请少量种子用户试用，针对脚本播报、推流稳定性和互动问答效果进行反馈收集。修复bug，调整AI语音或动画效果参数，完成第一阶段验收。

• 第二阶段：升级支持3D数字人 (预计2个月)

里程碑2.1 - 3D数字人形象制作接入（~1个月）：引入3D数字人生成与渲染引擎（例如Unity/Unreal或定制3D引擎）。实现从照片自动生成3D模型的流程⁴，或集成第三方高精度数字人模型库。解决3D人物的骨骼绑定和表情驱动，将语音驱动拓展到3D头部模型（利用Audio2Face等技术完成语音到表情的映射）。

里程碑2.2 - 3D场景和动作丰富化（~0.5个月）：开发虚拟场景切换、简单肢体动作/手势呈现等功能，使3D数字人不仅嘴巴动，还能有人物姿态变化。优化渲染效率，保证3D模式下实时推流性能。

里程碑2.3 - 项目验收与上线（~0.5个月）：整合2D与3D两套功能，在后台支持数字人形象的2D/3D切换。编写使用文档和培训材料。经过一轮全面稳定性测试和内容安全评估后，上线发布正式版本。

时间预估与资源

整个项目周期约为5~6个月。一阶段3个月交付可用的2D数字人直播MVP产品，二阶段2个月扩展3D能力。各阶段将投入以下资源：

- 人员投入：产品经理1人（全程统筹需求和验收）、后端工程师2人（分别负责AI模块和推流&后台模块）、前端工程师1人（负责编辑器和控制界面）、视觉设计/UI 1人、AI算法工程师1人（负责TTS和动画模型调优）、测试工程师1人。第二阶段增加1名3D技术美术/建模支持。
- 硬件&环境：开发测试使用高性能PC（GPU用于训练或推理加速唇动模型），云端部署预留服务器资源用于语音合成等服务。OBS软件及目标直播平台账号准备就绪。

里程碑验收标准

- 里程碑1.2结束时，2D数字人应能读出示例脚本并通过OBS在内部搭建的RTMP服务器观看到直播画面，视为核心链路打通。
- 里程碑1.4内测时，至少完成**10场**不同脚本的连续直播测试，收集到的主要问题均已修复，达到第一阶段验收标准所列各项指标。
- 里程碑2.2完成时，3D数字人形象外观质量和唇动同步效果需达到与2D相当的水平，渲染帧率≥30fps，无明显卡顿。
- 最终上线验收需通过安全审核（内容健康、无敏感违规）和压力测试（模拟高并发观看，系统稳定），并获得试用客户的认可反馈。

通过阶段划分和里程碑控制，确保项目在既定时间内稳步推进，各关键节点成果明确，为最终成功上线提供保障。

3. 产品方案设计

核心功能设计说明

结合PRD功能模块，这里进一步阐述各核心功能的方案设计细节：

- **数字人自动生成方案**：系统提供易于使用的人物生成向导。在用户提交照片时，后端调用人脸重建模型生成高精度的3D人头模型和对应贴图⁴；同时基于深度学习驱动的人脸关键点提取算法，自动完成骨骼绑定与权重设定，使生成的人脸能够支持后续表情动画²¹。对于半自动模式下的2D形象，系统内置多套美术绘制的角色立绘，可动态替换局部元素（如不同发型PNG图层），并通过插值合成实现某些中间变化，提升自定

义灵活性。所有数字人形象数据（2D的各图层或3D模型文件）将存储在数据库或文件存储中，并与用户账号关联。

- **AI语音合成方案**：采用业界领先的神经网络TTS模型，实现逼近真人的高逼真语音⁹。中文语料可使用科大讯飞的高品质语音合成服务或开源模型（如FastSpeech2+MB-iSTFT-VITS）做定制微调，确保发音贴合主播风格。为了满足情感表达需求，模型支持调整语速、音调和添加停顿、重读标签等。系统还允许用户录制少量语音样本，利用声音克隆技术训练专属音色，从而使数字人声音更独特（可选功能）。合成的音频在播放前做缓存，保障直播时不出现断续等待。
- **表情与动作驱动方案**：2D数字人采用音频驱动口型技术，即根据TTS生成的音频特征实时生成连续的口型帧动画。具体实现上，集成开源项目如Wav2Lip或其改进模型MuseTalk¹²对数字人原始头像进行唇部区域的合成，输出每帧对应的口型变化视频片段。相比传统逐帧手动画或Viseme插值，AI驱动方式省去手工制作，且能达到30fps实时效果²²。MuseTalk模型通过对音频编码与图像编码融合，使用轻量UNet推理，可在V100 GPU上实现帧率30以上²²，确保满足直播要求。对于其他表情和动作，2D数字人主要通过几种预设状态切换（如点头、微笑等简单动作）实现，这些状态可根据语义或情感分析触发，例如当讲到高兴内容时触发微笑表情。第二阶段的3D数字人将采用专业引擎驱动：通过NVIDIA Audio2Face或Unity的面部表情映射插件，将音素映射到3D模型的BlendShape权重上，实现与2D类似的嘴型同步。²³提到虚拟数字人系统包含语音生成和动画生成模块，实际我们的TTSA方案将语音和动画一体化驱动3D角色²⁴。
- **互动对话方案**：对于互动模式，设计上遵循ASR→NLP→TTS→动画的流水线架构¹¹。当观众以文字评论提问时，直接进入NLP意图识别步骤；当观众通过语音提问（如平台连麦或语音弹幕功能），则先调用自动语音识别(ASR)服务将语音转成文本²⁵。NLP部分如果只是FAQ匹配，则采用关键词+意图分类模型匹配问答库；如果启用开放AI对话，则将问题提交给大语言模型（如ChatGPT或本地部署的中文对话模型）生成回答¹¹。为确保回答准确，系统可以在LLM提示中加入本领域知识（例如产品参数、库存信息）作为上下文。生成的回答文本再交由TTS模块合成语音并驱动数字人回答。整套流程设计需要考虑实时性，各模块优化后使首个回应延迟控制在3秒左右¹⁴。在用户角度，将感觉数字人几乎是即时地听懂并作出了回答。



例如，百度的数字人直播平台在内容创作流程上经历了纯人工→半自动→全自动的迭代，不断引入AI以提高自动化程度²⁶。1.0人工阶段需要人工撰写脚本和手动应对问答；2.0半自动阶段引入AI辅助装修场景、创作脚本和托管部

分问答；3.0全自动阶段可由AI生成大部分内容并托管整个直播 27 8。本系统的方案也将提供手动与AI结合的模式：运营人员可选择让AI自动产出大部分内容，再人工审核调整，属于半自动创作；或完全信任AI，由其全权生成脚本和回答，实现高度自动化。通过分级的AI辅助，既保证内容质量又提升效率。

- **OBS推流集成设计：**由于不同平台推流机制各异，方案上采取**抽象推流接口**设计：
- 提供统一的“开始直播”控制按钮，内部根据配置选择具体的推流实现路径。对于抖音、小红书等无法直接推流的平台，前端会引导用户启动官方直播软件并启用OBS虚拟摄像头；系统后台监测到虚拟摄像已启用信号后，开始将画面输出到本地虚拟摄像设备，等同于摄像头输入，被官方软件采集 15。这个过程对于用户是透明的，但需要在产品使用文档中指导设置。
- 对于允许自定义RTMP的目标，如一些电商平台或海外TikTok等，系统直接使用推流库向多个RTMP地址推送。同一时间可以建立多路推流线程实现**一键多播**，也支持用户配置只推送选定的平台。
- 在技术实现上，集成**FFmpeg**命令或OBS的SDK库来抓取系统生成的画面帧和音频帧，然后编码为H.264/AAC流推送。鉴于实时性要求，每帧画面最好在16ms内处理完毕，我们会控制数字人动画生成和视频编码总耗时低于帧间隔。在网络层面加入断线重连机制，确保推流中断时自动重试。
- 此外，考虑到直播画面可能不仅有数字人，还有商品展示、字幕等叠加元素，我们可在OBS中预先设置这些素材源，数字人视频作为其中一个源进行合成推流。也可以让系统直接输出带有叠加内容的画面。
- **安全与审核设计：**产品方案特别关注内容合规，内置多层次的审核：
- 在脚本生成阶段，集成敏感词检测，对于AI自动生成的脚本内容先进行违禁词过滤和语义审查，必要时标记提醒用户修改。
- 在直播进行时，实时监控数字人将要说出的文本（因为我们控制TTS输入文本），可在最后一步拦截不当言论（例如LLM意外生成了不合适回答，则不予朗读）。同时监控弹幕问题，如果含有政治敏感/辱骂等，可选择不回答或礼貌拒绝。
- 系统管理员可配置违禁词列表和审核等级，通过后台干预紧急停止直播的接口，一旦发现问题可立刻中断推流。
- 用户数据安全方面，采取云端加密存储用户定制的形象和脚本，隐私信息不对外泄露。

用户体验流程图

为了直观展示用户使用本系统的体验流程，建议绘制整体**用户流程图**，分为“直播前配置”和“直播中互动”两部分：

直播前配置流程图：

用户登录 -> 创建/选择数字人形象 -> 准备直播内容

- > [选择模式：脚本直播 或 互动直播]
- > （脚本模式）编辑或生成脚本 -> 设置场景素材
- > （互动模式）配置问答库/知识库 -> 设置场景素材
- > 预览效果 -> 一键启动直播（进入推流）

上述流程图反映用户需要按顺序完成的步骤。比如用户如果选择纯脚本模式，可忽略问答库配置；若选择互动，则脚本可选简略提纲，由数字人自由发挥。预览效果是可选的校对环节。

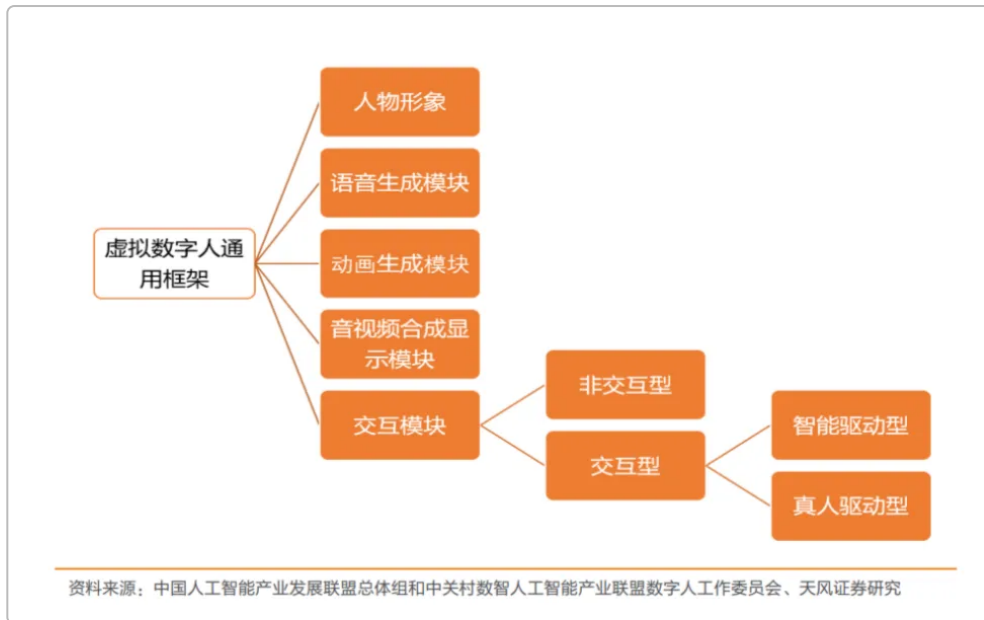
直播中互动流程图：

观众发送弹幕提问 -> 平台弹幕接口 -> [系统获取到问题]
 -> NLP判断：是常见问题？ -> 是：从问答库检索答案 -> 否：提交LLM生成答案
 -> 得到回答文本 -> 语音合成 -> 数字人播报回答（推送到直播间）
 -> 循环等待下一个问题 ...
 （脚本模式并行：若一段内容播报完且无新提问，自动进行下一段脚本）

这个流程图展现了系统在直播进行时的信息流转：从观众提问到数字人回答的闭环。¹¹ 所述的ASR->LLM->TTS->THG模块串联正对应于此。在设计流程图时，可将各模块标注清楚，例如“语音识别模块”“对话决策模块”等，使技术实现与用户视角对应。

以上流程图将有助于开发和产品人员对用户体验路径形成共识，确保各环节无遗漏且衔接顺畅。在实际设计UI时，也会据此流程引导用户逐步完成配置，不会出现跳步骤或迷失的情况。

4. 技术方案设计



整体技术架构：系统的技术架构可以分为前端应用、后端服务和外部集成三大部分。前端提供用户界面用于配置数字人和直播参数；后端包含数字人驱动核心逻辑和内容生成AI服务；外部集成包括OBS推流和第三方平台接口。整体架构如下图所示，典型的虚拟数字人系统由人物形象、语音生成、动画生成、音视频合成展示以及交互模块五部分组成²⁸。其中交互模块可选，如果启用则数字人分为交互型（进一步分智能驱动或真人驱动）和非交互型两类²⁹。本系统实现中，2D/3D数字人形象模块提供模型和素材支持，语音与动画模块构成“文本驱动语音和动作（TTSA）”引擎，由AI智能驱动；交互模块实现智能驱动型的实时对话功能，而无需真人驱动²⁴。最终通过推流模块将合成的音视频输出到各直播平台呈现给观众。

后端核心架构

后端架构采用模块化设计，各模块通过异步消息或调用管线串接，保证实时性能： - **内容管理子系统：**基于Web应用框架（如Node.js或Django）实现，负责存储和提供脚本文本、问答库、以及数字人形象配置等数据。提供API接

口供前端调用（例如获取某用户的数字人列表、保存新的脚本内容等），并在直播过程中将当前需要播报的文本发送给AI驱动模块。 - **AI驱动子系统**：这是系统的大脑，包含以下子模块：

文本分析/NLP模块：对输入的用户文本（脚本或弹幕）进行分析。如果是脚本段落，可能进行分句、插入适当停顿标签；如果是观众提问，进行意图分类和实体提取，调用问答或LLM服务拿到答复¹³。

语音合成模块：采用AI语音合成SDK或本地部署模型，将待播报文本转成语音数据（音频流）。优先使用本地模型以降低延迟，可选用开源实现或如讯飞在线TTS⁹。

动画生成模块：根据生成的语音，调用唇动驱动模型合成对应的数字人动画帧。2D场景下，这可能是通过Python子进程运行Wav2Lip模型得到短视频片段；3D场景下，则通过引擎插件直接驱动avatar模型的动画组件。

这些模块可串行工作：即当有一段文本需要播报时，NLP处理后送TTS，TTS输出音频后立即送动画模块开始处理。

而对于互动对话模式，不断循环等待新的输入触发上述流程。 - **直播控制模块**：负责与OBS等交互。实现方式可以通过OBS提供的WebSocket API来动态控制场景切换、开始停止虚拟摄像等操作。或者更简单地，系统输出一个本地**虚拟摄像头设备**（可用OpenCV+DirectShow等技术创建）供OBS采集¹⁵。在RTMP推流模式下，此模块直接使用FFmpeg库将接收到的动画帧与音频编码推流。直播控制模块还维护直播状态机：例如当前是播放脚本模式还是等待问答模式，并根据状态决定从内容管理取脚本还是等待NLP新问答，不断协调AI驱动模块工作。

模块之间通过事件队列或消息总线通信：例如，当TTS完成语音合成，会触发一个“AudioReady”事件，动画模块监听该事件后开始生成视频；又如当弹幕模块收到新问题，会触发“NLPRequest”交给NLP模块处理。这种架构使得各AI组件可以解耦独立替换升级，符合**多模块级联**的思想¹¹。

为保障性能，AI模型推理可使用GPU加速。如果部署环境支持，也可将TTS和唇动模型加载常驻内存（作为服务进程），避免重复加载模型耗时。对于LLM对话，可选云端API调用，以换取更高质量回复，但要做好网络超时和错误处理。

技术选型

结合项目需求和最新行业动态，选择合适的技术方案： - **AI语音合成**：优先考虑成熟的商用服务如科大讯飞的**在线语音合成**，其提供多情感音色选择，且在中文领域效果优秀，能满足实时性要求¹⁰。备选方案是本地部署基于Transformer或Diffusion的TTS模型，如微软的FastSpeech系列或VITS变体，经过高质量数据微调，可以生成接近真人的语音⁹。考虑到需支持多语言，项目中可集成多套模型或服务（讯飞适合中文，Google TTS或Azure TTS适合英文等）。 - **人脸动画驱动**：2D情况下选用**Wav2Lip**开源模型，该模型通过GAN训练对任意人物图像实现精准唇形合成，已被广泛验证效果³⁰。不过Wav2Lip可能分辨率有限（144p-240p嘴部），可尝试使用其改进版本

TalkLip或本项目提到的**MuseTalk**，后者在清晰度和实时性上更优¹²。3D情况下，优先利用**NVIDIA**

Audio2Face技术或游戏引擎的内置方案，这些方案通常通过预先绑定好的BlendShape来驱动3D面部表情，实时性强且避免训练复杂度。 - **自然语言处理**：问答匹配使用**Elasticsearch**或**Milvus**向量数据库结合嵌入模型，实现语义检索找答案，以提高对同义提问的覆盖率。开放式对话推荐接入OpenAI的GPT-4或国产大模型（如百度文心一言）作为后端¹¹，通过统一接口封装，未来可根据效果切换模型。鉴于实时对话要求高，可限制生成长度或选用压缩LLM（如ChatGLM等本地模型）以减少延迟。 - **OBS推流**：在系统实现上可以直接使用**OBS Studio**软件，通过启动OBS进程并操控其API完成推流；或者使用**FFmpeg**等工具库直接推流。考虑到OBS有完善的UI供用户调整推流参数，我们采取人机结合的方式：**推荐用户安装OBS**，系统通过文档指引用户获取推流密钥或使用虚拟摄像头接入¹⁵。对于有编程能力的平台（比如Twitch国外版），也可用RTMP直接推送，届时封装FFmpeg命令调用即可。 -

数据库：采用关系型数据库（如MySQL/PostgreSQL）存储系统结构化数据。库表设计主要包括：用户表（User），数字人表(Avatar)，脚本表(Script)，问答知识表(QnA)，直播记录表(StreamLog)等。其中Avatar表存储每个数字人的属性（形象文件路径、音色配置等）；Script表存储脚本段落文本及关联的Avatar和场景；QnATable存FAQ对；StreamLog记录每场直播的元数据（开始时间、时长、观看量等）和互动内容日志（可存JSON形式的问题->答案对列表）。此外，素材文件（头像图片、3D模型文件等）可存于云存储，由数据库保存路径引用。

系统示意与案例

为验证技术路线的可行性，我们参考目前业界和开源的类似方案。例如，开源项目**Open-LLM-VTuber**已经展示了本地实时语音对话驱动Live2D形象的能力³¹。该项目集成了语音识别、LLM和TTS，并通过Live2D实现了一个可在桌面实时交互的AI伴侣³²。这证明了我们的交互式架构在技术上是可行的。而商业产品**BocaLive**则将AI数字人用于电商直播，支持多平台推流和AI生成脚本，印证了我们的功能设计在商业场景的价值³⁵。

综上所述，本技术方案通过模块化架构和成熟AI技术选型，能够支撑一个AI语音驱动的数字人直播系统。从2D到3D的演进路径清晰，既利用现有技术快速实现，又为未来拓展打下基础。在实现过程中，我们会持续根据测试反馈优化模型效果（如让数字人讲话更自然、表情更丰富）以及系统性能，确保最终产品达到预期目标，为用户带来创新的数字人直播体验。

参考文献：

- 1. 百度数字人直播创编平台从人工到全自动的演进²⁶⁸
- 2. 中国电信分钟级全自动3D数字人生成技术报道⁴
- 3. 腾讯新闻：虚拟数字人系统框架与交互形式分析²³²⁴
- 4. 凤凰网科技：开源数字人实时对话Demo技术解析¹¹¹²
- 5. GitHub Open-LLM-VTuber项目简介³¹³²
- 6. BocaLive官方：AI数字人直播功能介绍³⁵
- 7. OBS推流官方教程对抖音平台的适配方案¹⁵
- 8. CSDN博客：AI数字人涉及的核心技术综述⁹¹³

128129130131132

BocaLive - AI Digital Human Live Streaming Software Platform
<https://www.bocalive.ai/>

10

421

中国电信发布新一代3D数字人：几张照片就能快速生成 超逼真_手机新浪网
<https://finance.sina.cn/tech/2023-11-11/detail-imzuhcep2481115.d.html?from=wap>

82627

百度数字人直播创编体验改版实战复盘！ | 人人都是产品经理
<https://www.woshipm.com/ai/6194966.html>

91325

AI数字人需要涉及多种技术_数字人开发所需技术-CSDN博客
https://blog.csdn.net/qq_15821487/article/details/130883723

1112142230

开源数字人实时对话：形象可自定义，支持语音输入，对话首包延迟可低至3s | 已上线阿里ModelScope魔搭社区_凤凰网
<https://i.ifeng.com/c/8dubLpCYERN>

151617181920

OBS直播推流技巧_OBS教程_OBS直播推流软件中文站 OBS官网版本分流下载 OBS插件免费下载
<https://www.obsproject.com.cn/obs/378.html>

23242829

硬核解析，一文看懂虚拟数字人的原理与机会_腾讯新闻
<https://news.qq.com/rain/a/20230618A059DR00>

3132

GitHub - Open-LLM-VTuber/Open-LLM-VTuber: Talk to any LLM with hands-free voice interaction, voice interruption, and Live2D taking face running locally across platforms
<https://github.com/Open-LLM-VTuber/Open-LLM-VTuber>