



实时语音对话与打断功能：OpenAI、Google 和 MiniCPM-o 的实现机制

简要概述

- **OpenAI ChatGPT语音交互：**OpenAI通过将大型语言模型（ChatGPT）与语音输入输出模块相结合，实现了语音对话功能。用户说话时，系统利用Whisper语音识别模型将语音转录成文本供ChatGPT理解，然后由ChatGPT生成回答文本，再通过新一代TTS（文本转语音）模型将回答朗读出来¹。该方案本质上还是串联管线：语音->文本->ChatGPT->文本->语音，各模块配合提供语音对话能力。ChatGPT的语音对话目前主要在移动App上实现，支持包括英语在内的多语言语音输入（依赖Whisper的多语种能力）和多种逼真的合成声音输出¹。用户需点击按钮开始讲话，整个过程接近实时但通常为半双工交互（一次只听或说）。若用户在ChatGPT朗读回答时开始说话，应用会停止播报并重新听取用户输入，实现一定程度的用户打断。不过，由于ChatGPT核心模型并非原生支持音频，其对语音的处理主要依赖前后的ASR/TTS模块，实时性能和打断灵活性相对有限。
- **Google Gemini Live (Project Astra)：**Google的Gemini系列模型从设计上原生支持多模态，能够直接理解和生成语音等多媒体信息²。其Project Astra原型展示了近乎实时的语音对话能力，将直播音频流输入AI模型，并以极低延迟产生回应³。Gemini 2.5模型可以直接“用音频思考和作答”，支持24种语言的语音输入输出，并能识别用户语音中的情感、口音等细节⁴。该系统具备全双工对话能力：对话过程中用户可随时打断AI讲话，AI会即时停止输出并倾听⁵。为实现这一点，Google采用了流式识别和超低延迟TTS技术，使语音对话像人与人交流一样顺畅⁶。“语音活动检测+语义判断”用于确定用户讲话何时结束或何时插话，以避免双方互相打断⁷。整体运行主要依托云端的Gemini模型和语音服务，通过Android设备等终端采集音视频流。Google还针对移动设备提供了不同等级的模型（如Gemini Pro/Flash）以权衡性能和延迟，并开放了Live API供开发者构建低延迟语音交互应用⁸。
- **MiniCPM-o 实时语音对话模型：**MiniCPM-o 2.6是清华大学等开源社区推出的仅8亿参数的端侧多模态大模型⁹。它在小模型中集成了视觉、语音和文本处理能力，支持图像、语音输入以及语音、文本输出，实现接近人类的低时延实时交互⁹。其架构采用端到端全模态设计：融合了已有的中型语音识别模型（如Whisper-medium）、语音合成模型（ChatTTS）和语言模型基座（Qwen 7B）等，通过联合训练使模型能够直接从语音waveform等输入推理并生成语音输出¹⁰。MiniCPM-o使用流式多模态机制来处理连续音频/视频流：将时间轴切分成小片段，在每个片段内并行处理新的音视频输入，同时判断是否该输出回复¹¹。这种时分复用并发方案确保模型能一边“听”一边“想”和“说”，极大降低交互延迟¹¹。当用户讲话尚未结束时，模型通过高层语义判断来预测句尾，从而几乎在用户停顿的瞬间就开始回答，减少等待¹²。如果对话中用户突然插话打断模型，MiniCPM-o也能立即响应，中断当前语音输出并转入聆听状态¹³。得益于端到端训练，模型还保留了语音信号中的情绪、语调信息，使生成语音更自然且可控风格：可通过提供参考声音（“语音系统提示”）让模型以指定音色和情感语调回答¹⁴。MiniCPM-o经过高度优化，可在平板电脑、手机等本地设备实时运行，实现真正的离线语音助手。在2025年CES展示中，它已在AI手机、车载终端等设备上流畅运行，被誉为“端侧GPT-4o级”技术方案^{9 15}。

OpenAI ChatGPT语音交互的技术实现

模型架构与多模态支持：ChatGPT本身是纯文本的大语言模型，并不直接支持语音输入输出。OpenAI采用多模块串联架构，将语音处理模块与ChatGPT文本对话模块结合¹。用户语音首先由Whisper模型进行语音识

别，Whisper将语音信号转成文本序列，再将该文本提交给ChatGPT生成回答。生成的回答文本再送入OpenAI自研的**文本转语音（TTS）**模型合成语音播放¹。这种架构下，各模块各司其职：Whisper负责语音->文本，ChatGPT负责对话逻辑，TTS负责文本->语音。ChatGPT无需改变内部结构即可通过这种管线“间接”支持语音，因此无需对核心模型进行多模态训练。但代价是语音中的非文字信息（如情感、语气）在转换为文字时可能丢失，模型无法利用。

低延迟音频处理：为提升对话实时性，OpenAI对语音管线进行了优化。Whisper在移动端App中以**流式方式**工作，边听边出部分转写结果，减少等待时间（一些开发者实测Whisper小模型每秒可实时转录，延迟仅几十毫秒¹⁶）。ChatGPT API本身支持**流式输出**，所以App一拿到ChatGPT部分回答文本就立刻开始合成和播放，无需等待整句完成。这种**流水线并行处理**（识别-生成-合成同步进行）实现了接近实时的响应。此外，系统运用了**端点检测和语音活动检测（VAD）**：当用户停止讲话出现静默时，Whisper或独立VAD模块判断讲话结束，立即触发ChatGPT生成回答，减少不必要的等待。Whisper具备多语言识别能力，在安静环境下英文识别延迟可低至百毫秒级，使ChatGPT的回应能够在用户说完后很快开始。在语音输出方面，OpenAI的TTS模型经过优化，**首字输出延迟很短**，合成语速接近实时讲话速度¹⁶。总体而言，虽然ChatGPT语音交互走的是串联管线，但通过流式处理和并行化，大幅降低了感知延迟，实现了自然的对话节奏。

打断机制：ChatGPT语音交互最初采用**半双工**交互模式：用户按住麦克风按钮说完话-> ChatGPT回应完毕-> 再轮到用户讲话。这种设计避免双方同时说话，但缺乏人性化的灵活性。为改善体验，OpenAI加入了一定的**打断支持**：当ChatGPT正在语音回答时，若用户开始说话（或再次按下话筒），系统**立即中断**语音播报并切换回听取用户输入状态。这个过程利用了**并行监听技术**——即使ChatGPT在播报，应用仍持续监测麦克风输入的能量变化或唤醒词。一旦检测到用户声音，便触发中断。由于Whisper可以持续转录长达30秒的音频，理论上应用可以在ChatGPT讲话时后台运行Whisper监听用户插话。从用户角度看，就是可以随时打断AI的长篇回答，提出新问题。在处理打断时，ChatGPT生成端也会收到停止信号，中断当前回答生成。这种机制依赖可靠的**VAD**来区分用户有意的插话 vs. 背景噪音。相比Google和MiniCPM-o的全双工方案，ChatGPT的打断能力属于应用层面的权宜之计，但已经让对话更加**自主可控**。值得注意的是，ChatGPT尚未提供像人那样的“预测对方话语结束”能力，仍主要依靠静音阈值来判断用户是否说完，因此偶尔会有等待过长或过早打断的问题。

语音识别与合成模块：OpenAI选择了业界知名的**Whisper**模型作为ASR模块¹。Whisper是端到端模型，支持多语言高准确率识别，这保证了ChatGPT能听懂不同语言和口音的用户输入（目前主要支持包括英语、法语、汉语在内的几十种语言）。Whisper模型有不同规模版本，为兼顾速度和准确率，移动App可能使用中等大小的模型并在服务端运行，借助云端算力提升实时性。语音合成方面，OpenAI研发了**高度逼真的TTS技术**¹。他们与专业配音演员合作，打造了5种独特的AI声音，每种声音通过几秒真人语音克隆训练而来¹。这套TTS模型能生成几乎乱真的人声，包括停顿、韵律和情感，使ChatGPT听起来更自然。此外，生成语音的速度经过优化，短文本几乎瞬时合成，长文本也能边生成边播放，不拖慢对话节奏。可以说，Whisper保证了“听得懂”，TTS保证了“说得像”，共同支撑了ChatGPT的语音对话体验。

部署与优化：ChatGPT语音对话目前主要通过**移动端应用+云服务模式**提供。语音数据和文字请求都会发送至OpenAI云服务器处理，这意味着需要网络连接。将计算放在云端可以利用强大的GPU/TPU资源运行大型模型（如GPT-4）和高精度的语音模型，以获得最佳效果。移动App本地只需完成音频录制、回放等轻量工作。不过，为进一步降低延迟，OpenAI也可能在终端上执行部分任务：例如在手机端利用神经芯片跑一个小型VAD或短语识别模型，快速响应打断等交互，再将完整音频上传云端获取最终结果。一些第三方实现已经展示了在PC本地运行Whisper和GPT-3等模型进行语音对话的可行性¹⁶。但OpenAI官方尚未发布轻量本地版本的大模型。可以预见，未来随着模型蒸馏和设备算力提升，ChatGPT语音助手也有望推出**精简版**，在保护隐私和低延迟方面提供更优的本地体验。

Google Gemini Live (Project Astra) 的语音对话技术

多模态原生架构：Google的Gemini模型系列从底层架构上就是**多模态的**。不同于OpenAI将语音作为外挂模块，Gemini在训练时就让模型**直接吃音频、图像等输入**²。官方指出Gemini是“从零开始构建的原生多模态模型，能跨文本、图像、音频、视频理解和生成内容”²。在Gemini 2.5中，语音对话被视作核心能力之一，

模型可以 **直接以语音形式推理和回答**¹⁷。推测其技术路线与Google早前研究AudioPaLM类似：采用统一的**Token表示**，将文字和语音都编码成序列输入Transformer。例如，用户的语音可能通过一个音频编码器（类似Whisper或Conformer模型）转成一串声学tokens，这些tokens与文本token共存在模型上下文中，让模型同时利用语音内容和声音特征信息进行理解。模型在输出时，可以直接生成表示语音的token序列，再经解码器合成音频。这样，整个对话无需中途转换为纯文本语义，语音中的说话人情绪、语气等**副语言信息**也能被模型感知并体现在回答中¹⁸。此外，Gemini通过多模态训练还能让语音和视觉、文本信息互相协同：例如看着视频画面进行讲解，用语音描述看到的内容¹⁹²⁰。总之，Gemini/Astra是一个**端到端多模态大模型**，在同一个Transformer框架内实现了听觉、视觉和语言的融合，这为实时对话打下了基础。

实时语音流处理：Project Astra之所以能实现**低至几乎无延迟**的语音对话²¹，归功于Google在流式AI上的深厚经验。首先，Google采用**Streaming ASR**技术：其语音识别采用改良的RNN-T或Transformer流模型，可以在用户讲话过程中持续输出部分转写结果（Assistant的识别甚至做到边听边在字幕上滚动文字）。Gemini模型能一边获取转写文本，一边开始思考回应。在I/O 2025展示中，Astra被证明能对用户摄像头看到和麦克风听到的实时流做出瞬时反应²¹。其次，**流式推理**：Gemini可能使用了分块推理或增量编解码技术，使模型不必等完整一句话结束才给出反应。例如，对话系统可以在检测到用户问题的大致意图后，就提前检索或准备部分回答内容。当用户话音一落，系统几乎同步地开始语音回答。Google官方提到Astra利用了Gemini的**Flash推理**能力：这是Gemini的一种高效模式，可在保持大部分质量的同时显著加快生成速度⁶。最后，Google通过大规模工程优化（模型剪枝、并行计算、专用硬件）将延迟进一步压低。据TechCrunch报道，Astra实现了**音视频流输入到答案输出几乎无时滞**，仿佛AI就在现场对话²¹。这种性能在很大程度上依赖云端强算力和模型优化，同时Gemini Live移动端应用也做了本地**快速响应**设计，如动画过渡、partial response提示等，提升主观流畅度。整体而言，Google的方案把**端到端延迟压缩**到人类对话可接受范围内（几百毫秒量级），真正实现了“**实时**”语音交流。

支持打断与全双工对话：在自然对话中，双方可以随时插话、打断。Project Astra/Gemini特别强调了“**不打断的顺畅交流**”和“**随时响应**”⁷。这意味着系统具备**全双工**能力：AI边说话时仍在“听”用户。一旦用户开口，AI将**立即停下**，切换为聆听模式，从而不会与用户声音重叠⁵。实现这一点需要多方面配合：1) **持续监听**：Gemini Live在AI说话的同时，通过设备麦克风+音频处理模块持续监测环境声音。2) **语音活动检测(VAD)**：如果探测到明显的用户语音信号（音量、频谱特征超过背景噪音阈值），立即判定为用户想说话。3) **即时中断TTS**：系统能够暂停或停止正在播放的合成语音。这可能通过将TTS输出分成小音频块流式播放，一旦检测打断就停止发送后续块，实现毫秒级停播。4) **上下文管理**：当用户半途插话时，AI需丢弃尚未播报完的剩余回答，并根据新输入调整后续对话走向。Gemini模型强大的上下文记忆可以确保即使被打断，也能衔接对话，不会“忘记”先前内容。5) **预测结束点**：更高级地，系统可能训练了**语音端点预测模型**（类似电话系统里的Turn Taking模型），根据用户语音内容和语调来预判TA是不是快说完了。如果模型确信用户的问题已完整表达，它就会迅速开始回答，哪怕用户话音刚落或还有轻微犹豫。这避免了人为加长的停顿。Google在学术上有过相关研究（如双向Attention的Transducer模型，可使用一点点后文预测流结束）。在Astra系统中，他们提到通过**模型的高级语义判断用户语音结束时机**¹²，这佐证了端点预测的重要性。6) **忽略背景干扰**：为了防止误打断，Gemini训练时学会了**区分直接对话语音和环境噪声/他人说话**²²。例如，用户旁边有人在聊天，Astra不会把那当作指令，也不会乱插话回应。因此在没有被呼叫时，AI会安静等待（正如官方描述“明白什么时候不说话”¹⁸）。综合这些机制，Gemini Live实现了像人一样知进退的对话风格：不抢话，也不让用户等，支持多人场景下只对主人回应。这种全双工打断能力远超传统语音助手，带来了更自然的交互体验。

语音识别与合成模块：Google在语音AI领域深耕多年，其ASR和TTS技术业界领先。Gemini/Astra很可能融合了Google现有的**最佳ASR模型**。当前Google有自研的**Conformer**架构流式识别，以及**Whisper**类似的大模型。Gemini 2.5本身据称**原生支持语音识别**，可以处理多达24种语言输入⁴。它应能直接输出转写文本或语义表示，无需借助外部服务。此外，Astra强调了**情感识别**能力，能够察觉用户语气中的愤怒、愉快等情绪⁸。这可能在ASR层面附带了**情感分类**，或者Gemini模型对音色做了联合训练，使其隐变量中带有人声情感信息。TTS方面，Google在I/O发布了**可控风格TTS**技术，允许通过提示调整语调、语速、情绪等²³²⁴。Gemini 2.5提供**Pro**和**Flash**两种TTS生成模式：²⁵ Pro侧重高保真音质，Flash兼顾速度成本。开发者可以根据应用选择，用更快模型实现低延迟响应。TTS输出支持**多语言、多说话人**，甚至可以让AI一人分饰两角对话²⁶。这些能力在Gemini Live上体现在：AI声音更加生动（重音、节奏贴近真人），可按需调整风格（比

如“用更愉快的语气回答”之类的指令），并支持多语言混合对话而无需切换模型²⁷。由于Gemini能理解用户的情感和语气，它的TTS回应也会相应做出匹配。例如用户生气时，AI可能用平缓安抚的语气回答（官方称之为感情对话能力²⁸）。在底层实现上，Google很可能采用**两阶段TTS**：Gemini先输出带语调标记的文本或中间表示，再交由专门的神经声码器（如WaveRNN/Parallel WaveNet或最新的SoundStorm模型）生成波形。无论实现细节如何，Google的全栈语音AI技术保证了Astra系统在听清用户话、准确理解、多彩发声各方面都达到业界顶尖水准。

运行平台与优化：Gemini Live主要作为**云端服务**运行，通过手机、眼镜等终端获取输入输出。Gemini 2.x模型参数庞大（据传Pro版上千万参数），不可能在手机本地推理，所以用户语音和摄像头画面会**加密上传至**Google数据中心，由Gemini和相关语音服务实时处理，然后将结果下发设备播放。为了降低用户设备的负担和延迟，Google在设备侧也有一些优化措施：例如**部分本地AI计算**——Android设备可能使用内置DSP或神经引擎跑一个**低功耗唤醒词检测**或短语识别，让设备在本地判断用户是否在叫“Hey Google”或已经在讲话，从而及时打开麦克风、开始上传数据。还有报道指出，Google正与三星等合作研发**Astra智能眼镜**^{29 30}。眼镜这类设备对时延更敏感，预计会引入**边缘计算**：即在手机端（配套眼镜的手机）部署**小型Gemini模型**或推理加速器，承担一部分推理任务（尤其是视频帧的初步处理、压缩编码发送）。事实上，Google DeepMind也推出了**Gemma**系列轻量模型，参数量较小，可以离线运行某些功能。未来Gemini或Astra可能通过模型裁剪/蒸馏，推出能在高端手机NPU上跑的简化版，实现部分离线交互。不过在当前，完全的GPT-4级别对话还是需要云计算支持。Google通过其全球分布的服务器和高速网络，实现了Astra服务在不同地区的**低延迟接入**。例如，与用户最近的数据中心承担推理，加上传输优化，可以把语音往返时间降到很低，让用户感觉不到明显网络延时。总之，Google采用**云+端协同**策略：云端提供强大的多模态AI能力，端侧处理传感和部分智能辅助，以实现流畅实时的语音对话体验。

MiniCPM-o 模型的语音对话实现

端到端多模态架构：MiniCPM-o 2.6最大的特点是将**多模态能力集成进一个小模型**。它以一个8B参数的Transformer同时处理视觉、听觉和语言信息⁹。开发团队将几个现有模型“融合”成一体并联合训练：¹⁰显示MiniCPM-o以SigLip（可能是视觉-语音特征模型）、Whisper-medium ASR、ChatTTS语音解码器和Qwen2.5-7B语言模型为基础，构造成统一架构。与传统串联方式不同，MiniCPM-o选用**端到端全模态训练**：语音不先转文本、视频不先抽帧特征，而是让模型直接学习从**原始音频、视频输入到目标输出**的映射³¹。为了做到这一点，研究者对离线的编码器/解码器进行了**流式改造**，使其适应在线数据流³²。比如，用改进的ViT模型来编码视觉帧序列，用改进的Whisper来编码音频流，得到的表示不再只是文字转写，还包含声音音色等信息³³。这些多模态表示接入共同的Transformer层，与文本token一起计算，从而模型内部实现了不同模态语义的**对齐融合**³⁴。在解码端，MiniCPM-o有一个**自回归语音解码模块**负责生成语音：即模型输出不是普通文本字符，而是一系列可直接合成语音的单位（可能是声码器codec的codes或者带音高时长标注的梅尔频谱参数）。通过这种架构，MiniCPM-o可以将**音频/视频中无法文字描述的细节**也携带到最后输出。例如，用户语气中的戏谑、不满，这些信息虽然纯文字看不出，但模型因为端到端训练过，能捕捉到并反映在回答语音的情绪上^{32 35}。这一架构也让MiniCPM-o具备**声音克隆能力**：给模型一个几秒的说话人语音作为条件，它就能用这个声音说出回答^{36 14}。总的来说，MiniCPM-o通过创新的架构设计，在小小8B参数内打通了多模态语义鸿沟，实现了高度集成的“视觉-听觉-语言”理解与生成，为实时对话打下了基础。

低延迟的流式处理：在端侧设备上实现实时多模态对话，需要精巧的流式处理机制。MiniCPM-o提出了“**全模态流式**”方案¹¹：模型将来自不同模态的并行数据流按时间片交织处理，实现类人实时感知与响应。具体而言，它将时间轴划分为循环的短片段（比如每0.5秒一个片段）³⁷。在每个时间片内：模型一方面对这一时段内采集的视觉帧和音频信号做**增量编码**，更新对环境的认识；另一方面判断是否该在此刻输出语音回复。如果用户正在说话，模型持续编码他的语音；如果用户停顿或提问结束，模型的**决策单元**会在某个时间片触发**生成回复**。生成阶段同样以流式进行：MiniCPM-o不会一次性憋出长段回答再朗读，而是一边生成一边播放，使回复更及时。在用户听感上，MiniCPM-o往往能在用户话音刚落的不到半秒内开始回答，几乎无缝衔接。这得益于**它语音结束预测机制**：据介绍，模型利用其语言理解判断用户提问是否完成，从而避免等不必要的额外静默¹²。例如用户问“明天北京天气怎么样”，模型可能在听到“怎么样”之前就已根据上下文预测这是句完整的问题，提前在后台准备回答。当确认用户声音结束瞬间，模型立即说出答案。另一方面，如果用户发言冗长

或含糊，模型也不会贸然打断，而是持续倾听，通过语义和语音信号综合判断何时接话最合适。如此设计确保了既不抢话也不滞后。另外，MiniCPM-o针对端侧设备有限的显存和算力，也采用了一些高效算法：如超高模态像素密度技术，将每帧图像编码所需token从常规的2000+压缩到640³⁸（通过视觉特征稀疏表示等方法），使模型可低成本“看”更多视频帧；以及RAG（检索增强）机制管理长时输入³⁹——对于超出上下文窗口的长视频音频，模型将已处理部分摘要存储，必要时检索，避免一次性占用大量内存。这类似于人为的短期记忆+长期记忆结合，让模型在流式处理时既快又“记得住”。再加上8B参数模型本身计算相对更快，这些都使得MiniCPM-o在边缘设备上达到接近GPT-4级别的实时性能⁴⁰。

打断与全双工交互：作为面向人机自然交互设计的模型，MiniCPM-o支持随时打断和连续对话。官方描述中，它能够在被打断后及时作出反应，并以恰当的情绪语调继续回复¹³。其流式架构天生具备全双工能力：因为模型每隔几十或上百毫秒就在轮询音频输入，所以即使正在输出语音，也不会错过新的用户声音信号。一旦检测到用户插话，MiniCPM-o会立刻中止当前语音解码，转而解析新的用户语音输入。这个过程中延迟极低，因为无需切换模型或暂停线程，模型本身就是在交替进行听和说。可以想象其工作方式类似电话里的对讲机被改造造成开放麦克风：AI说话时仍在“听”着，只是降低输出优先级，一旦听到对方声音就迅速停下来。为了做到这一点，MiniCPM-o结合VAD和语义判断：VAD负责检测音频能量变化，语义判断则让模型确认这是不是用户在对模型说话（而非背景响声）。此外，MiniCPM-o通过系统训练让模型学会礼貌等待：如果用户声音未停，它不会贸然插嘴。这种对人类对话礼仪的学习，使AI的打断行为显得更自然。值得一提的是，MiniCPM-o还能依据对话状态调整情绪语调，包括在被打断时表达适当的反应（比如语气上表现出停顿或一丝错愕，然后转为聆听）。当用户说完重新让AI说时，模型又可以用不同情绪继续——比如之前场景是激动的，AI在被打断后回应时可能语气有所缓和或疑问，以匹配新的对话走向。这些细节都源于模型对对话上下文和语音信号的综合把握。通过全双工流式机制，MiniCPM-o在端侧实现了媲美专业语音助手的打断处理，而且因为一切在本地完成，延迟和依赖均更可控。

语音识别与合成：MiniCPM-o在语音理解和生成上均达到了开源模型的SOTA水准⁴¹。在语音识别方面，它支持中英双语口语识别，准确率甚至超过一些专门的ASR模型（如阿里Qwen2-Audio-7B）⁴¹。这说明其内部可能集成了Whisper模型的权重作为初始，并在对话数据上进一步微调，使之更适应交互场景的口语、噪声、口音等。由于端到端训练，MiniCPM-o识别时不只输出文本，还将音色、情绪等“听”到的信息保存在隐层，为后续生成服务³²。语音合成方面，MiniCPM-o采用自回归解码输出语音。推测它利用了像EnCodec这样的神经压缩码，将波形表示成离散token序列，并令模型学习预测这些token。这样8B的小模型也能生成高质量声音，因为最后的声音细节复原交由外部小型解码器完成（类似于Stable Diffusion先输出图像latent再解码成高清图）。事实证明MiniCPM-o的语音自然度和多样性极佳：不仅音质接近真人，还能在情感和音色上自由变化。官方演示包括模仿知名人物声音（特朗普和麦当劳叔叔）³⁶、在不同情境下用惊慌或愤怒的语气说话¹⁴等，体现了一人千面的声音塑造能力。这通过“语音系统提示”实现：开发者可以提供一段参考音频给模型，作为生成时的条件，就像指定了角色的声音。模型会解析该音频的说话人特征，并在合成输出时应用相同的特征¹⁴。此外，通过自然语言也能控制风格，比如在prompt里说明“请用幽默的语调回答”，模型会相应调整语音表现。这种高度可控的TTS能力，在开源模型里尚属少见，得益于MiniCPM-o打通了语音和语言的表示。最后，其语音合成速度在端上经过优化：精简的8B模型和高效声码器使得延迟极低。据报道，MiniCPM-o语音对答已经可以达到人讲话的实时速度，与GPT-4o云端语音接口相当⁹。综合而言，MiniCPM-o通过融合ASR和TTS模块，并用对话数据对其进行强化训练，实现了小模型的大能力：听得懂、说得出，还说得像。

端侧部署与模型优化：MiniCPM-o的另一大亮点是可在本地设备实时运行。8B参数大小约占数十MB显存，已在高端平板、手机上跑通（例如iPad Pro的M系列芯片可以推理8B模型）。面壁智能团队对模型进行了多方面优化：采用更高效的算子实现、INT4/INT8量化权重、剪枝蒸馏等，使模型在移动SoC/NPU上达到可接受的速度。他们在CES 2025上展示了MiniCPM-o在AI手机、车载系统、智能音箱等多种边缘设备上的实时应用¹⁵。端侧运行带来诸多优势：首先是低延迟，因为省去了云通讯时间；其次是隐私友好，用户语音不必上传服务器；再次是离线可用，在无网络时AI语音助手也能工作。当然，为在小设备上达到高性能，模型做了取舍：8B参数虽精炼，但与GPT-4数百亿参数相比仍有差距。一些复杂推理任务上MiniCPM-o可能不及云端大模型。不过，其设计目标更多是满足日常对话和多模态助手功能，因此在这些场景下通过专项优化达到了接近大模型的效果⁴⁰。未来，随着芯片性能提升和算法进步，我们可能看到更大的端侧模型（如20B、50B）出现。但MiniCPM-o已经证明，通过密度优化和创新架构，端侧设备完全可以承载高度智能的实时语音助手。这预示着大模型下沉终端的趋势：让每个人的手机里都有一个不依赖云的贴身AI对话助手。

延伸阅读与参考资料

- **OpenAI 官方博客**: 《ChatGPT can now see, hear, and speak》 – 介绍了ChatGPT新增语音交互功能的实现 (Whisper用于语音识别、新的TTS模型用于生成语音等) [1](#)。
- **Google AI 博客**: 《Advanced audio dialog and generation with Gemini 2.5》 – 详述了Gemini 2.5在语音对话和生成方面的的新能力，包括原生音频对话、低延迟、多情感语音合成等 [6](#) [18](#)。
- **Google DeepMind 项目页**: 《Project Astra》 – 展示了Astra原型的功能亮点，如**自然交互**（多语言流畅语音对话、无打断）、**情境感知**（忽略背景谈话）等 [42](#) [7](#)。
- **TechCrunch 科技新闻**: 《Project Astra comes to Google Search, Gemini, and developers》 – 报道了Google在I/O 2025宣布将Astra技术应用于搜索和Gemini等产品的消息，强调了Astra实现的**低延迟、多模态实时能力** [43](#) [8](#)。
- **Android 官方文章**: 《Explore new camera and screen-sharing in Gemini Live》 – 介绍Gemini Live的新特性，其中明确提到**用户可在任意时刻打断Gemini对话**，强调了对话的自然流畅 [5](#)。
- **搜狐科技**: 《仅8个月就把GPT-4o带到了端侧，面壁智能拿到了什么秘籍？》 – 深入报道了MiniCPM-o 2.6模型的技术细节，包括端到端全模态架构、流式机制、打断与情感语音等实现，以及其在端侧设备上的表现 [9](#) [12](#)。
- **甲子光年（微信号）** 对上述MiniCPM-o的报道 – 从研发故事角度剖析了面壁智能如何在短时间内打造端侧GPT-4o级模型，内容与搜狐报道类似，可一并参考。
- 其他参考资料：如OpenBMB社区提供的MiniCPM-o开源说明、知乎专栏“语音大模型概述”等，也对比分析了各家语音大模型的能力和实现原理，适合进一步阅读研究。 [10](#) [36](#)

[1](#) ChatGPT can now see, hear, and speak | OpenAI

<https://openai.com/index/chatgpt-can-now-see-hear-and-speak/>

[2](#) [6](#) [17](#) [18](#) [19](#) [20](#) [23](#) [24](#) [25](#) [26](#) [28](#) Gemini 2.5's native audio capabilities

<https://blog.google/technology/google-deepmind/gemini-2-5-native-audio/>

[3](#) [8](#) [21](#) [29](#) [30](#) [43](#) Project Astra comes to Google Search, Gemini, and developers | TechCrunch

<https://techcrunch.com/2025/05/20/project-astra-comes-to-google-search-gemini-and-developers/>

[4](#) [7](#) [22](#) [27](#) [42](#) Project Astra - Google DeepMind

<https://deepmind.google/models/project-astra/>

[5](#) Gemini Live: Use Camera & Screen Sharing on Android | Android

<https://www.android.com/articles/gemini-on-android/>

[9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [31](#) [32](#) [33](#) [34](#) [35](#) [36](#) [37](#) [38](#) [39](#) [40](#) [41](#) 仅8个月就把GPT-4o带到了端侧，面壁智能拿到了什么秘籍？ | 甲子光年_MiniCPM-o_模型_视频

https://www.sohu.com/a/849808088_100016644

[16](#) Show HN: Real-time AI Voice Chat at ~500ms Latency | Hacker News

<https://news.ycombinator.com/item?id=43899028>