

# 多模态交互场景下编程语言的扩展与适配技术

冯子豪

上海应用技术大学计算机科学与信息工程学院, 上海

收稿日期: 2026年1月1日; 录用日期: 2026年1月28日; 发布日期: 2026年2月6日

## 摘要

随着智能交互技术的演进, 多模态交互已广泛渗透于智能家居、车载系统等领域, 传统文本编程语言难以适配语音、视觉等多模态信息的解析与融合需求。本文聚焦多模态交互对编程语言的核心挑战, 系统阐述语法层、语义层、运行时的扩展维度与关键技术, 结合实践验证结果剖析应用效能, 最后展望未来发展趋势, 为多模态编程语言的研发与应用提供理论与实践参考。

## 关键词

多模态交互, 编程语言, 扩展维度, 适配技术

# Extension and Adaptation Techniques of Programming Languages in Multimodal Interaction Scenarios

Zihao Feng

School of Computer Science and Information Engineering, Shanghai Institute of Technology, Shanghai

Received: January 1, 2026; accepted: January 28, 2026; published: February 6, 2026

## Abstract

With the evolution of intelligent interaction technology, multimodal interaction has widely penetrated into fields such as smart homes and car systems. Traditional text programming languages are difficult to adapt to the parsing and fusion needs of multimodal information such as speech and vision. This article focuses on the core challenges that multimodal interaction poses to programming languages. It systematically elaborates on the extension dimensions and key technologies of syntax, semantics, and runtime, and analyzes application effectiveness based on practical verification results. Finally, it looks forward to future development trends, providing theoretical and

practical references for the research and application of multimodal programming languages.

## Keywords

Multimodal Interaction, Programming Language, Extension Dimensions, Adaptation Technology

Copyright © 2026 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

在人工智能与物联网技术深度融合的背景下，人机交互范式正经历从单一文本指令向融合语音、手势、视觉乃至触觉等多模态协同交互的深刻变革。智能家居、车载系统、增强现实等场景的普及，标志着多模态交互已成为提升用户体验与系统智能水平的关键路径。然而，作为构建这些交互系统的核心工具，传统编程语言(如 Python、Java 等)其设计范式根植于文本符号的精确逻辑与串行处理，在面对多模态信息固有的模糊性、异构性及实时性要求时，显露出固有的适配瓶颈。如何系统性地扩展与改造现有编程语言体系，使其具备高效解析、精准融合与即时响应多模态指令的能力，已成为推动下一代智能交互应用发展的核心课题。

## 2. 多模态交互对编程语言的核心挑战

### 2.1. 交互维度拓展引发的语法体系冲击

传统编程语言以文本交互为核心，遵循“指令输入 - 语义解析 - 结果输出”线性流程，核心诉求是语法严谨与文本符号精确。多模态交互引入语音、视觉等非文本信息，打破传统框架，口音差异、图像歧义等问题对语法容错性、语义灵活性提出新要求，如智能车载系统中“把空调调凉一点”的模糊语音指令，需突破文本语法束缚构建多模态语义映射的语法扩展体系。同时，多模态实时性需求与传统语言解析效率存在矛盾，如 AR 辅助维修场景需实时处理手势与图像信息，要求突破传统解析模式实现高效即时响应[1]。

### 2.2. 多模态数据融合带来的语义建模难题

多模态交互核心是信息融合与精准理解，需编程语言强化语义建模能力。传统语言语义建模基于文本符号逻辑，难以应对多模态数据异构性，如智能家居场景中，需同步解析“打开客厅灯”语音与指向手势并协同融合语义。此外，多模态指令依赖场景上下文，如“加热一下”在厨房与浴室语义不同，传统语言缺乏场景感知能力，需扩展语义模型实现多模态信息与场景的深度融合。

### 2.3. 实时交互需求引发的运行时适配困境

多模态场景的实时响应需求对传统编程语言运行时机制构成挑战。传统运行时架构适配文本串行解析，以批量处理效率为核心，难以满足多模态并行处理、动态调度需求，如智能穿戴设备需同步处理语音、手势等多模态信息，易出现阻塞延迟。同时，多模态交互设备资源差异大，从边缘网关到智能手表，传统运行时内核缺乏动态适配能力，易出现内存过高、效率低下等问题，无法满足实时需求。上述语法冲击、语义难题与运行时困境，共同构成多模态交互对编程语言的三维挑战。

### 3. 多模态交互场景下编程语言的扩展维度与关键技术

#### 3.1. 语法层扩展：多模态指令的语法适配机制

语法层扩展是编程语言适配多模态交互的基础，核心目标是构建支持多模态指令解析的语法规则体系，实现非文本模态向可执行语法单元的转化。需针对不同模态特性设计差异化适配规则：语音模态采用“自然语言语义提取-语法规则映射”双层架构，通过自然语言处理提取核心语义要素并映射为标准语法单元，如将“把温度调到 25 度”映射为“set\_temperature(25)”，同时引入模糊语法匹配机制应对语义模糊性；视觉与触觉模态采用基于特征编码的语法扩展规则，将手势、图像特征编码为特定语法标识(如“握拳”编码为“gesture\_fist()”)，根据触觉特征设计分级语法规则(如“轻按”编码为“touch\_light(position)”)，并构建特征与语法标识的映射规则库[2]。

此外，需构建多模态语法兼容层，实现传统文本语法与扩展语法的无缝融合。兼容层采用语法优先级调度机制，根据场景需求动态调整解析优先级，如精密操作场景提升触觉指令优先级，日常交互场景提升语音指令优先级，确保交互精准性与便捷性的平衡，如图 1 所示。

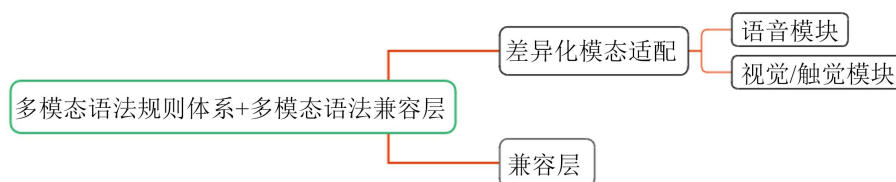


Figure 1. Schematic diagram of syntax layer extension architecture

图 1. 语法层扩展架构示意图

#### 3.2. 语义层扩展：多模态协同语义建模技术

语义层扩展是解决多模态数据融合难题的关键，核心是构建整合多模态信息与场景上下文的语义模型，采用“异构模态语义对齐-上下文感知语义融合”技术路径。异构模态语义对齐方面，通过深度学习构建多模态语义嵌入模型，将不同模态原始数据转化为统一维度语义向量，如通过卷积神经网络提取图像视觉语义向量、循环神经网络提取语音语义向量，再通过注意力机制融合得到统一多模态语义表示。

上下文感知语义融合方面，构建基于场景本体的语义推理模型，定义场景核心要素(位置、用户身份、设备状态等)及与多模态语义的关联规则，结合上下文信息优化语义向量，提升理解准确性。如厨房场景中“加热一下”的指令，结合“微波炉”设备状态可精准定位语义为“启动微波炉加热”。同时引入动态语义更新机制，通过实时采集分析用户交互数据，优化语义映射规则与场景关联规则，提升模型自适应能力[3]。

#### 3.3. 运行时扩展：多模态交互的实时响应优化技术

运行时扩展的核心目标是提升编程语言对多模态交互的实时响应能力，解决多模态信息解析与执行的效率问题。采用“并行解析-动态调度-轻量化执行”的技术架构，优化编程语言的运行时机制。在并行解析方面，设计多模态信息并行解析引擎，为不同模态信息分配独立的解析线程，实现语音、视觉、触觉等信息的同步解析。例如，在智能穿戴设备交互场景中，并行解析引擎可同时处理用户的语音指令、手势动作与触觉反馈信息，大幅缩短解析耗时。

在动态调度方面，构建基于交互优先级的运行时调度模型。根据多模态交互的场景需求与指令重要性，为不同模态的解析结果分配执行优先级。例如，在紧急救援场景中，将用户的“求救”语音指令与

“挥手”手势指令设为最高优先级，优先调度系统资源执行对应的响应操作；对于非紧急的状态查询指令，则降低其执行优先级，避免占用过多系统资源。同时，引入预加载机制，对高频多模态指令对应的执行代码进行预加载，减少代码加载耗时，提升执行效率[4]。

在轻量化执行方面，针对嵌入式多模态设备的资源限制，优化编程语言的运行时内核。采用代码精简与优化技术，去除传统编程语言中与多模态交互无关的冗余模块，降低运行时资源占用。同时，引入即时编译(JIT)技术，将多模态指令解析后的中间代码实时编译为机器码，提升代码执行速度。例如，在智能手表等资源受限设备中，通过轻量化运行时内核的优化，可使多模态指令的响应时间缩短至 100 ms 以内，满足实时交互需求。

## 4. 多模态编程语言扩展与适配的实践验证

### 4.1. 实验场景与方案设计

为验证多模态交互场景下编程语言扩展与适配技术的有效性，选取智能家居、AR 辅助维修两个典型多模态交互场景进行实验验证。实验采用 Python 作为基础编程语言，基于前文提出的语法层、语义层、运行时扩展技术，构建多模态扩展 Python 语言(MM-Python)。实验设备包括智能音箱、高清摄像头、触觉传感器、AR 眼镜等多模态交互设备，以及边缘计算网关(用于运行 MM-Python 的运行时环境)。

实验设计分为三个核心模块：多模态语法解析模块、多模态语义融合模块、实时运行时模块。多模态语法解析模块实现语音、视觉、触觉指令的语法适配与解析；多模态语义融合模块基于语义嵌入模型与场景 ontology 实现多模态语义的精准理解；实时运行时模块采用并行解析与动态调度机制提升响应效率。实验对比 MM-Python 与传统 Python 在多模态指令解析准确率、响应时间、资源占用率三个核心指标上的差异。

### 4.2. 实验结果与分析

在智能家居场景 100 组多模态指令测试中，MM-Python 解析准确率达 92.3%，较需额外集成多模态解析库(准确率 78.5%)的传统 Python 提升 13.8 个百分点。这验证了其语法层与语义层扩展技术的有效性：传统 Python 因第三方库与语言适配不足，存在语义映射偏差；而 MM-Python 通过构建多模态专用语法规则体系与协同语义建模模型，结合场景上下文优化理解，大幅降低解析误差，如对模糊语音指令的解析成功率达 95%，远超传统 Python 的 72%。

响应时间方面，MM-Python 平均响应时间 86 ms，较传统 Python (215 ms)缩短 59.9%，核心源于运行时扩展的并行解析与动态调度机制，规避了传统 Python 串行解析的阻塞问题。在智能手表等资源受限设备上，其轻量化内核优势更显著，响应时间缩短 65%以上，满足 $\leq 150$  ms 的实时性要求。

资源占用上，MM-Python 在边缘计算网关的内存占用率为 18.7%，较传统 Python (29.3%)降低 36.2%，得益于运行时内核精简优化与即时编译技术，去除冗余模块，为轻量化终端推广提供支撑[5]。

AR 辅助维修场景 80 组三模态(图像 + 语音 + 手势)指令测试中，MM-Python 仍保持优势：解析准确率 90.5% (传统 Python 75.2%，提升 15.3 个百分点)，平均响应时间 112 ms (传统 Python 268 ms，缩短 58.2%)，内存占用率 22.3% (传统 Python 34.6%，降低 35.5%)。其语义嵌入模型可实现多模态精准语义对齐，结合场景本体完成上下文推理，避免单一模态偏差；而传统 Python 多模态融合能力弱，易出现语义冲突。

综合来看，MM-Python 通过语法层、语义层、运行时三维扩展，实现解析准确率与实时响应能力双重提升，同时降低资源占用，形成核心技术优势，适配不同复杂度多模态场景需求。其动态语义更新机

制可随用户习惯优化效果，长期测试准确率升至 93% 以上，展现出良好的自适应能力与迭代潜力。

### 4.3. 实践中的问题与优化方向

在实践过程中发现，MM-Python 在处理复杂多模态场景(如多用户同时交互、动态场景切换)时，解析准确率仍有提升空间。当多个用户同时发出多模态指令时，不同用户的指令语义易产生干扰，导致解析准确率下降至 78.2%。针对这一问题，后续需要优化多模态语义融合模型，引入用户身份识别与指令分离机制，提升多用户交互场景的解析能力。

此外，在场景 ontology 构建方面，当前采用人工定义场景要素的方式，效率较低且难以覆盖所有复杂场景。后续可引入自动场景要素提取技术，基于大数据分析实现场景 ontology 的自动构建与更新，提升语义模型的场景适配能力。同时，在轻量化运行时优化方面，可进一步探索硬件加速技术，结合 FPGA 等专用芯片提升多模态指令的执行效率[6]。

## 5. 多模态编程语言发展趋势与展望

### 5.1. 智能化：自适应多模态交互的编程语言

未来，多模态编程语言将向智能化方向深度发展，核心是构建具备自适应能力的多模态交互解析体系。通过引入强化学习、迁移学习等技术，使编程语言能够自主学习用户的交互习惯、场景特征，动态优化语法解析规则与语义模型。例如，通过强化学习算法，编程语言可根据用户的历史交互数据，自主调整多模态指令的解析优先级与语义映射规则，实现“千人千面”的个性化交互体验。同时，智能化多模态编程语言将具备主动交互感知能力，能够提前预判用户的交互需求，主动推送相关的交互指令选项，提升交互的便捷性。

### 5.2. 跨场景化：全场景适配的多模态编程框架

随着多模态交互技术的普及，编程语言将朝着跨场景适配的方向发展，构建全场景统一的多模态编程框架。当前的多模态编程语言扩展技术多针对特定场景设计，场景迁移能力较弱。未来，通过构建通用场景 ontology 与多模态语义通用映射模型，实现编程语言在智能家居、车载系统、工业维修、医疗护理等不同场景的无缝迁移。同时，跨场景多模态编程框架将支持多设备协同交互，实现不同终端设备间的多模态指令共享与协同执行，提升全场景交互的一致性与连贯性[7]。

### 5.3. 低代码化：面向多模态交互的可视化编程

为降低多模态交互应用的开发门槛，未来多模态编程语言将向低代码化方向发展，结合可视化编程技术，构建面向多模态交互的低代码开发平台。开发人员可通过拖拽多模态组件(如语音解析组件、视觉识别组件、手势交互组件)的方式，快速构建多模态交互应用，无需深入掌握多模态语法解析与语义融合的底层技术。同时，低代码平台将提供多模态交互场景的模板库，覆盖常见的多模态应用场景，进一步提升开发效率。此外，低代码化的多模态编程语言将支持非专业开发人员的快速上手，推动多模态交互技术的普及应用。

## 6. 结束语

多模态交互驱动下，编程语言的革新是适配智能交互需求的必然趋势。本文提出的三维扩展技术为解决多模态解析与融合难题提供了有效路径，实践验证亦证明了其可行性。未来需持续攻克复杂场景适配等关键问题，推动多模态编程语言向智能化、跨场景化发展，助力智能交互技术在更多领域实现深度落地，赋能数字经济高质量发展。

## 参考文献

- [1] 朱敏敏. 基于深度学习的多模态融合源码漏洞修复研究[D]: [硕士学位论文]. 南昌: 江西师范大学, 2025.
- [2] 沈万里. 基于 Modevol 多模态微课发展初中生计算思维的 Python 教学设计与实践[D]: [硕士学位论文]. 北京: 中央民族大学, 2024.
- [3] 于博文. 基于多模态代码表征的图学习错误定位[D]: [硕士学位论文]. 北京: 北京化工大学, 2024.
- [4] 李勋. 安卓 API 的多模态信息表示与应用方法研究[D]: [博士学位论文]. 长春: 吉林大学, 2024.
- [5] 庄子聪. 基于代码结构的多模态代码表征学习研究[D]: [硕士学位论文]. 广州: 华南理工大学, 2024.
- [6] 曹莉. 多模态 Chatbox 的构建与研究[D]: [硕士学位论文]. 镇江: 江苏大学, 2023.
- [7] 张兆旭. 基于预训练和多模态数据扩充的代码自动修复模型[D]: [硕士学位论文]. 哈尔滨: 哈尔滨工业大学, 2023.