



上海交通大学硕士学位论文

基于面部捕捉、语音、头部运动的在线实时数
字人驱动

申请学位：工学硕士

学科/专业：计算机科学与技术

院 系：计算机学院

2026 年 1 月 8 日

**A Dissertation Submitted to
Shanghai Jiao Tong University for the Degree of Master**

**FACECAPGES: REAL-TIME FRAME-BY-FRAME
GESTURE GENERATION FROM AUDIO, FACIAL
CAPTURE, AND HEAD POSE**

School of Computer Science
Shanghai Jiao Tong University
Shanghai, P.R. China
January 8th, 2026

上海交通大学

学位论文原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下，独立进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不包含任何其他个人或集体已经发表或撰写过的作品成果。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本人完全知晓本声明的法律后果由本人承担。

学位论文作者签名：

日期： 年 月 日

上海交通大学

学位论文使用授权书

本人同意学校保留并向国家有关部门或机构送交论文的复印件和电子版，允许论文被查阅和借阅。

本学位论文属于：

- 公开论文
 - 内部论文，保密 1年 / 2年 / 3年，过保密期后适用本授权书。
 - 秘密论文，保密 ____ 年（不超过 10 年），过保密期后适用本授权书。
 - 机密论文，保密 ____ 年（不超过 20 年），过保密期后适用本授权书。
- （请在以上方框内选择打“√”）

学位论文作者签名：

指导教师签名：

日期： 年 月 日 日期： 年 月 日

摘要

手势是人类交流中用于补充语义细节与传递情绪的重要非语言信号。在虚拟人/数字形象的实时控制中，传统动作捕捉往往依赖穿戴式设备，虽然精度高，但存在成本高、使用门槛高、便携性差等问题；相机动捕虽降低了硬件负担，却仍要求用户在镜头前持续表演手势，在远程交流、直播等长时间场景中会带来空间占用与体力消耗，并与键鼠操作产生冲突。近年来语音驱动手势生成技术为低门槛驱动提供了可能，但现有主流方法通常假设可获取完整语句或未来上下文，难以直接应用于用户逐字语音输入的在线实时交互环境，因此需要一种低延迟、可部署的实时手势生成方案。

为解决上述问题，本论文提出 FaceCapGes：一种面向实时数字人驱动的帧级手势生成方法，仅使用可在线采集的三类信号：语音、面部表情与头部姿态，在不依赖未来输入的条件下逐帧生成上半身 3D 骨骼动作，从而让用户无需实际做出手势即可驱动虚拟形象获得自然的随语表达。本文的主要研究内容包括：

(1) 提出 FaceCapGes 在线实时多模态手势生成框架，给出从信号采集、帧级推理到虚拟人驱动渲染的整体流程设计，并在严格因果约束下明确其可部署的实时生成机制。

(2) 在级联多模态架构中引入头部姿态作为新的实时输入模态，设计并实现轻量姿态编码与弱耦合融合策略，为手势提供节奏前瞻与空间锚定线索，从而增强动作节奏与空间指向一致性。

(3) 提出并实现在线实时手势生成的滑动窗口自回归训练策略，通过片段切割与窗口展开实现严格因果的逐帧生成，并结合片段级监督与对抗训练目标，提升实时生成的稳定性与时间连续性，缓解自回归漂移与抖动问题。

(4) 开展用户主观实验、客观指标测量与实时性能测试，并与代表性方法进行对比评估。实验结果表明，在严格因果约束下，本文方法 FaceCapGes 在真实性上表现良好，与现有扩散模型相当；同时在语调变化较快与头部朝向变化明显的片段中，能够生成更平滑且指向一致的动作，并具备满足实时交互场景需求的推理性能。

关键词：协同语音手势生成，数字人驱动，面部捕捉，多模态学习

Abstract

Gestures are essential nonverbal signals in human communication, enriching semantic details and conveying affective states. In real-time control of virtual humans or digital avatars, conventional motion capture systems typically rely on wearable devices. Although accurate, such systems are costly, inconvenient to use, and lack portability. Camera-based capture reduces hardware requirements but still forces users to continuously perform gestures in front of a camera, which occupies physical space, causes fatigue during long-term scenarios such as remote communication and live streaming, and conflicts with keyboard–mouse interaction. Recently, speech-driven gesture generation has emerged as a low-barrier alternative; however, most existing approaches assume access to complete utterances or future context, making them difficult to deploy in online interactive settings where speech is streamed word by word. Therefore, a low-latency and deployable real-time gesture generation solution is required.

To address these challenges, this thesis proposes FaceCapGes, a frame-level gesture generation method for real-time avatar driving. FaceCapGes uses only three types of online-available signals—speech, facial expressions, and head pose—to generate upper-body 3D skeletal motions frame by frame without relying on future inputs, enabling users to drive avatars with natural co-speech gestures without physically performing them. The main contributions of this thesis are summarized as follows:

(1) A real-time multimodal gesture generation framework, FaceCapGes, is proposed, providing an overall pipeline design from signal acquisition and frame-level inference to avatar driving and rendering, and clarifying a deployable real-time generation mechanism under strict causality constraints.

(2) Head pose is introduced as a new real-time input modality in a cascaded multimodal architecture. A lightweight pose encoder and a weakly coupled fusion strategy are designed and implemented to provide rhythmic anticipation and spatial anchoring cues, thereby enhancing gesture timing and spatial pointing consistency.

(3) A sliding-window autoregressive training strategy for online real-time gesture generation is proposed and implemented. By segmenting motion sequences and unfolding them with a sliding window, strictly causal frame-wise generation is achieved. Combined with

segment-level supervision and adversarial objectives, the proposed strategy improves temporal stability and continuity, alleviating autoregressive drift and jitter.

(4) Subjective user studies, objective metric evaluation, and real-time performance tests are conducted, with comparative assessments against representative methods. The results demonstrate that, under strict causality constraints, FaceCapGes achieves strong realism comparable to existing diffusion-based models; moreover, in segments with rapid prosodic variations and noticeable head orientation changes, it generates smoother motions with more consistent spatial pointing, while maintaining inference efficiency sufficient for real-time interactive applications.

Key words: co-speech gesture generation, virtual avatar driving, face-capture, multimodal learning

目 录

第 1 章 引言	1
1.1 研究背景和意义.....	1
1.2 手势生成质量评估点.....	2
1.3 国内外研究现状.....	3
1.3.1 规则驱动阶段	3
1.3.2 数据驱动阶段	3
1.3.3 多模态扩展阶段	4
1.3.4 国内外研究现状总结	4
1.4 在线实时手势生成与离线生成的差异.....	5
1.4.1 面向 AI 虚拟形象的整句手势生成.....	5
1.4.2 面向用户交互的在线实时手势生成	5
1.5 本文研究目标与研究内容.....	6
1.6 论文组织架构.....	7
第 2 章 多模态在线实时手势生成框架	8
2.1 手势的定义.....	8
2.1.1 Kendon 连续体	8
2.1.2 随语手势的分类	9
2.1.3 头部手势的分类	10
2.2 身体姿态的参数化表示.....	10
2.3 面部表情的定义与参数化表示.....	12
2.4 系统设计思路.....	13
2.4.1 手势生成目标的选择	13
2.4.2 引入头部姿态的动机与贡献	14
2.5 在线实时手势生成系统整体框架.....	16
2.5.1 信号采集与系统配置层	16
2.5.2 手势生成模型层 (FaceCapGes)	17
2.5.3 渲染驱动层	18

2.6 本章小结.....	18
第3章 融合头部姿态的多模态级联手势生成模型架构	19
3.1 问题定义.....	19
3.1.1 任务描述	19
3.1.2 输入与输出模态	20
3.2 级联架构设计的继承.....	20
3.2.1 级联架构的原理与理论背景	20
3.2.2 说话人 ID 分支移除.....	21
3.2.3 输入模态的解码	21
3.2.4 身体姿态的解码	22
3.3 头部姿态模态的引入.....	23
3.3.1 输入特征的定义	23
3.3.2 特征获取方法	23
3.3.3 级联架构中的位置	23
3.4 模型整体结构.....	24
3.5 本章小结.....	24
第4章 基于滑动窗口的实时手势生成自回归训练	26
4.1 单步因果 LSTM 预测器	26
4.1.1 LSTM 基本概念与状态传递机制.....	26
4.1.2 双时间尺度记忆结构	27
4.1.3 窗口内预测过程的形式化表达	28
4.2 训练片段切割.....	28
4.2.1 固定长度片段定义	28
4.2.2 重叠切割与样本覆盖	29
4.3 片段内部的滑动窗口展开策略.....	29
4.3.1 前置动作帧与预热阶段	29
4.3.2 滑动窗口展开与逐帧自回归生成	30
4.3.3 拼接生成序列与片段级输出	30
4.4 监督损失.....	31
4.4.1 片段级生成序列与损失计算范围	31
4.4.2 总体优化目标	31

4.4.3 姿态重构与时序平滑损失	31
4.4.4 损失权重设置	32
4.5 对抗训练.....	32
4.5.1 基于拼接片段的片段级判别	32
4.5.2 对抗损失定义	33
4.5.3 交替优化策略	33
4.6 本章小结.....	33
第 5 章 实验结果与分析	35
5.1 训练配置.....	35
5.2 实验配置.....	36
5.2.1 实验对比模型	36
5.2.2 跨模型评估设置	36
5.3 用户评估.....	37
5.3.1 用户评估系统与实验配置	37
5.3.2 结果与分析	40
5.4 定性分析.....	42
5.4.1 生成动作平滑性	42
5.4.2 头部朝向与手势空间指向一致性	43
5.5 客观评估指标与实现细节.....	45
5.5.1 Fréchet 手势距离 (Fréchet Gesture Distance)	45
5.5.2 语义相关动作召回率 (Semantic Relevance Gesture Recall)	47
5.5.3 节奏对齐度 (Beat Alignment)	47
5.5.4 L1 范数	48
5.5.5 评估区域设定与公平性说明	49
5.6 定量评估结果.....	49
5.7 消融实验分析.....	50
5.8 性能评估.....	51
5.8.1 单帧推理性能	51
5.8.2 端到端计算链路延迟	52
5.8.3 系统更新率	53
5.8.4 帧率设定的可扩展性	53

5.9 本章小结.....	54
第 6 章 结论	55
6.1 本文工作总结.....	55
6.2 未来工作展望.....	56
6.2.1 高层语义信息	56
6.2.2 面向未来趋势的预测性训练目标	56
参考文献.....	58
附录 A 手势生成对比视频	63
附录 B 代码与实现资源	64
学术论文和科研成果目录.....	65

插 图

图 2.1 FACS 中闭眼动作单元的定义	13
图 2.2 BlendShape 在闭眼形变上的线性插值效果	13
图 2.3 系统整体架构与数据流示意图	16
图 3.1 头部姿态编码器结构示意图	24
图 3.2 FaceCapGes 模型整体结构	25
图 5.1 BEAT 数据集的骨架拓扑结构与驱动范围	35
图 5.2 用户评估工具实机界面	39
图 5.3 用户评估总体主观排名结果	41
图 5.4 平衡拉丁方设置下的主观排名结果	41
图 5.5 生成动作效果对比	43
图 5.6 头部朝向与手势指向一致性的动作对比	44

表 格

表 1.1 离线与在线手势生成任务在约束与可利用信息上的对比.....	5
表 2.1 不同旋转表示方式的空间连续性与使用示例.....	12
表 3.1 输入输出模态符号与维度.....	20
表 5.1 对比模型的输入输出模态.....	37
表 5.2 定量评估结果.....	50
表 5.3 消融实验结果.....	50
表 5.4 推理速度评估结果.....	52

第1章 引言

1.1 研究背景和意义

近年来，随着元宇宙、虚拟社交与直播等领域的相关技术日趋成熟，用户已能够使用任意外观的虚拟人作为交互载体，在虚拟空间中与异地用户进行交流。虚拟人3D模型的姿态由其内部骨架关节的旋转参数（如欧拉角、四元数等）定义，最终通过蒙皮渲染技术完成可视化。得益于自动骨骼绑定技术，骨骼动画的生成可消除不同3D模型间的骨架拓扑差异，实现跨模型复用。

在虚拟人交互中，穿戴式动作捕捉设备是实时驱动手势的传统方案，贴合肢体的标记点可将肢体运动实时转换为相同的骨骼动画，提供直观、准确的操控。尽管其精度较高，但对于大多数用户而言，此类设备存在功能用途单一、硬件成本昂贵、便携性差等问题，限制了使用频率。因此，在当前的虚拟人交互应用中，仅少数专业用户会使用此类设备，而多数用户在设备限制下无法简单控制虚拟人肢体，导致了两种用户体验之间的不一致。

针对普通用户对低门槛虚拟人交互的需求，基于相机的动作捕捉技术^[1-2]成为主流替代方案。该技术无需额外硬件，仅通过手机或电脑的内置相机即可实时捕捉用户动作，用于转化为骨骼动画。然而，该方案仍存在三种局限：一是这种方法要求用户在空中面向相机做出手势，过程中难以同步操作键盘、鼠标，造成操作冲突；二是手势活动范围受相机视野限制，在自然使用距离通过手机或电脑的内置相机拍摄用户，将严重限制手势的捕捉范围；三是持续的手势动作会产生体力消耗，在直播等长时间使用场景中，用户疲劳问题将变得显著。

因此，我们提出一种新的需求：一种无需用户实际做出手势，仅通过用户的实时语音、面部捕捉与头部姿态，实时生成与语义和情感相匹配的手势骨骼动画。该方法旨在降低使用门槛，并解决操作冲突与体力消耗问题。

然而，现有研究尚未提供成熟的解决方案满足上述需求。首先，当前多数手势生成方法依赖完整的语音或文本输入。由于用户的语音输入是逐字进行的，计算机需要等待用户的未来输入才能解析当前的语义，造成动画生成的延迟。其次，现有研究未对用户的头部姿态模态做深入研究。头部姿态对手势的节奏与朝向具有明确的关系，且可以通过常规相机实时捕捉，但利用该模态来增强生成手势的自然度的相关研究尚不充分。

为此，本文提出了一种新颖的实时手势生成模型。该模型以帧为单位，输入语音、面部表情与头部姿态数据，并逐帧输出对应的骨骼动画。本文首次在实时手势生成中，将头部姿态作为一种新的模态引入，利用其与自然手势节奏与朝向的高度相关性，结合语音和面部信息共同提升生成动作的自然度。我们采用级联多模态架构与自回归训练来融合这些模态，以学习其联合表征，从而在严格的实时约束下增强手势的表现力。

1.2 手势生成质量评估点

语音驱动手势生成的目标并非仅复现关节轨迹，而是在交互语境中生成“看起来像人说话时自然会出现”的动作表现。

生成手势通常应满足以下直观属性：

1. 动作应具有自然的运动学规律与视觉连贯性，避免抖动、漂移与非人体的突然加速；
2. 手势的节奏与语音韵律应在时间上协调一致，使动作的起势、峰值与回落与重音或语调变化相呼应；
3. 生成结果应具有一定的变化性与表达张力，避免长时间静止或重复单一模式；
4. 在具备语义表达能力的任务中，动作还应与语义内容或情绪倾向保持一致；
5. 此外，面向实时交互系统的研究还必须兼顾低延迟与稳定性，使动作能够在连续输入流下平滑输出并保持可部署性。

基于上述目标，现有研究形成了较为体系化的评价维度，可概括为以下几类：

1. **自然性（Naturalness）** ——衡量生成动作在运动平滑性、速度变化及能量分布上的合理性，常采用 FGD (Fréchet Gesture Distance)^[3]、运动速度统计、或主观“自然度”评分等指标。
2. **同步性（Synchronization）** ——评估手势在时间上与语音重音或韵律事件的对齐程度，常用 BA (Beat Alignment)^[4]等方法，以及基于重读检测的主观同步性评价。
3. **多样性（Diversity）** ——衡量模型在不同语音输入下生成的动作变化程度，通常以轨迹分布的方差、速度曲线差异或 L1 范数等指标度量，以防止模型陷入单一模式或过度平滑。
4. **语义相关性（Semantic Relevance）** ——反映生成动作与语义关键词或情绪类别的一致性，可通过 SRGR (Semantic Relevance to Gesture Ratio)^[5] 等指标或人工

标注语义标签对齐评估。

5. **实时性与稳定性 (Latency & Robustness)** ——在面向交互系统的研究中，还需评估帧级推理延迟与输出平滑性，以确保动作流连续且系统响应及时。

1.3 国内外研究现状

近年来，语音驱动手势生成经历了从规则设计到数据驱动模型、再到多模态扩展与实时生成的持续演变。这一过程不仅体现了算法架构的更新，也反映了研究目标与应用场景的变化：从基于语言规则的行为映射，到学习语音—动作关系的深度生成模型，再到面向交互的多模态实时系统。

1.3.1 规则驱动阶段

早期的手势生成系统主要依赖语言学规则与专家知识构建^[6-9]。这类方法通过语义分类或韵律规则将语音片段映射为预定义的手势模板（如指示、肯定、节奏性动作），并以有限的动作库组合出手势序列。它们可在虚拟代理或机器人中实现基于语音的同步动作。然而，手势词典与语法规则的人工设计成本较高，难以覆盖自然语音中的多样变化，导致生成结果缺乏自然性与个体差异。

1.3.2 数据驱动阶段

随着大规模语音与动作配对数据的出现，研究者开始采用统计学习和深度神经网络模型学习语音一手势映射关系。在此阶段，语音通常作为唯一输入模态，模型通过长短记忆网络（Long Short-Term Memory, LSTM）、多层感知机（Multilayer Perceptron, MLP）等结构预测连续手势序列。典型代表如 CaMN 模型^[5]，其基于 BEAT 数据集^[5]训练级联网络，将 LSTM、全连接网络与 GAN 结构相结合，实现从语音到动作的端到端预测。

然而，该类模型多使用欧拉角或离散旋转参数作为手势表示，生成结果容易出现抖动与不连续。后续工作引入更平滑的表示方式，如 Rot6D^[10-13]或 Axis-Angle^[14]，显著提升了动作流畅性。与此同时，为解决语音与手势间的多对多映射问题，研究者引入了向量量化变分自编码器（Vector Quantized Variational AutoEncoder, VQ-VAE）^[11,15]与扩散模型^[14,16-19]，在保持自然性的同时提升了生成多样性与表现力。

尽管这些方法在客观指标与视觉效果上优于传统模型，但通常依赖完整语句级上下文。在用户语音下的流式逐字输入场景中，为获取未来上下文以进行语义判别与韵律对齐，需引入缓冲机制，因此即使推理较快的模型^[14]，整体延迟也因上下文缓冲

造成端到端的显著延迟。

1.3.3 多模态扩展阶段

为进一步提升动作表现力与语音理解能力，部分研究引入视觉模态或语言语义特征。例如，CaMN^[5]在语音输入的基础上融合面部捕捉信息以增强表现；EMAGE^[11]与DiffSHEG^[14]同时生成手势与面部动作；DiffTED^[18]实现了端到端的视频合成。

这些多模态生成方法在提升虚拟智能体的自然感与沉浸感方面表现优异，但其任务假设仍基于整句输入，因此主要用于AI虚拟形象生成或离线内容创作场景，而非实时用户交互。

1.3.4 国内外研究现状总结

本节关注近年来代表性生成范式及其在线化需求；其中部分方法在离线设定下效果突出，但本文受严格因果与低延迟约束未予实现，仅作为后续工作参考与对照基线。

技术趋势 近年来的手势生成研究在建模策略上呈现两类趋势：

1. **扩散模型（Diffusion-based generation）**：扩散模型在手势生成中常带来较高的动作自然度与多样性^[14,16-18,20]。现有工作多以固定长度的语音/文本片段作为条件，在迭代去噪采样中生成整段手势序列，因而在直接迁移到在线场景时会引入片段缓冲的延迟。另一方面，也已有研究开始探索离线实时或任意长度生成的采样与拼接策略（如基于扩展/外延采样的设计），表明扩散式框架具备向流式化演进的潜力^[14]。
2. **语义增强方向（Semantic-aware generation）**：为覆盖iconic、metaphoric等更依赖语义的手势，引入语义表征作为额外条件^[20-21]，以提升动作与语义的一致性与表现力。在线场景下，该方向的主要挑战通常来自语义信号的获得方式：现有工作使用整句/整段文本嵌入，基于更充分的上下文以稳定语义对齐，因而需要通过有限前瞻或缓冲策略来权衡延迟与语义质量。

总体而言，上述方法多在离线或半离线设定下达到最佳效果；若面向严格低延迟的在线实时生成，需配合流式条件建模、有限前瞻与快速采样等机制。本次将作为后续工作的技术储备与对照基线。

1.4 在线实时手势生成与离线生成的差异

现有手势生成研究可从面向用户的在线实时生成、与面向 AI 的离线生成分为两类，两者的差异在表 1.1 进行了对比：

表 1.1 离线与在线手势生成任务在约束与可利用信息上的对比

Table 1.1 Comparison of Constraints and Available Information between Offline and Online Gesture Generation

对比维度	AI 虚拟形象生成	用户虚拟人实时生成
输入信息	完整句级语音或文本（可使用未来信息）	实时语音流，仅使用过去与当前帧
输出目标	整句手势序列（离线生成）	连续流式手势（逐帧生成）
时间约束	不必要实时	帧级实时性（<50 ms 推理延迟）
评价重点	整体语义一致性与美学自然度	瞬时同步性、动作平滑与交互稳定性
应用场景	离线动画、内容合成、AI 虚拟直播	实时虚拟人、视频会议、用户虚拟直播
语音模态	作为输入或由文本生成（TTS 输出）	作为实时输入特征（语音流）
手部手势	生成目标（输出）	生成目标（输出）
面部表情	通常为生成目标（输出）	可通过设备实时采集，作为输入辅助推理
头部姿态	通常为生成目标（输出）	可实时采集并作为输入特征，用于同步推理

1.4.1 面向 AI 虚拟形象的整句手势生成

这一类研究的目标是让 AI 驱动的文本对话系统的虚拟形象具备手势表现力。

模型可以一次性生成下一句语音或文本，因此可以访问完整的未来信息，包括整句音频、文本和语义上下文。

典型方法通过编码完整句子的节奏与语义，预测整段动作轨迹，以最大化动作与语义的一致性和整体流畅性。

这类方法适合 AI 驱动的系统或离线生成的应用场景，如合成视频。

1.4.2 面向用户交互的在线实时手势生成

本文所聚焦的目标类型属于此类。

在用户实时说话的过程中，系统需根据当前语音流（以及可选的面部表情与头部姿态）即时生成同步手势。

此任务具有严格的实时性约束与因果性限制：模型在每一时刻只能使用当前及过去的信息，而不能访问未来语音或文本内容。

该任务更接近实时交互系统，而非内容生成系统。

1.5 本文研究目标与研究内容

本文研究的目标是设计一种能够在实时条件下运行的语音驱动手势生成模型，使用户无需动作捕捉设备或特定硬件，仅通过语音输入与相机前的面部表情、头部动作即可驱动虚拟人的上肢与头部动作。与以往主要面向离线生成或预先制作型虚拟形象生成的研究不同，本文关注的任务场景为用户实时交互，因此系统必须在严格因果的条件下运行，即仅利用当前与过去的输入帧信息进行动作预测，避免依赖未来语音或文本内容。此外，实时运行对推理延迟与计算效率也提出了更高要求。

为满足上述需求，本文进一步引入用户的头部姿态作为辅助输入模态，为手势生成过程提供额外的非语言信号。头部姿态能反映注意方向与交互焦点，并在语音节奏发生变化时为动作的时序组织提供参考。通过将语音、面部与头部信号联合输入模型，系统能够在实时条件下获得更丰富的上下文线索，从而支撑连续动作的稳定生成。

本文以 CaMN 模型^[5]为基础进行扩展。CaMN 原为离线级联结构，其输入包括语音与面部捕捉特征，输出包含手部上肢动作与头部姿态。本文将其输入机制改写为逐帧输入的流式推理形式，并在此基础上引入头部姿态特征分析模块，将头部姿态作为独立通道输入至级联网络的后级层，以实现语音、面部与头部信号的联合驱动。最终，系统能够在实时语音流输入条件下逐帧生成骨骼动画输出，满足实时交互场景对因果性与低延迟的要求。

因此，本文的研究内容汇总如下：

1. 多模态在线实时数字人手势生成系统设计：构建从设备平台采集语音、面部捕捉与头部姿态等输入模态，并以逐帧流式推理方式生成实时上半身动作的端到端系统流程，提出可部署的在线数字人驱动系统。
2. 融合头部姿态的多模态级联手势生成模型架构：将头部旋转姿态作为新的实时输入模态，引入适用于在线场景的特征表示与编码器设计，并与语音、面部表情特征进行联合建模与融合，形成继承 CaMN 级联解码策略的多模态动作生成结构。
3. 基于滑动窗口的自回归训练与在线推理一致性策略：针对严格因果的实时生成

约束，提出片段切割与滑动窗口展开的自回归训练流程，结合单向时序解码器与历史动作缓冲机制，使模型在仅依赖当前及历史信息的条件下生成连续稳定的动作序列。

4. 实验与评估平台实现：搭建统一的数据处理、推理与渲染流程，支持不同模型在相同输入条件下进行动作生成与公平对比，执行客观评价指标与主观用户实验的完整评估流程。

1.6 论文组织架构

本文章节安排如下：

第一章陈述在线实时数字人驱动场景下，语音驱动手势生成的研究背景与意义，概述生成质量的主要评价维度，总结国内外研究现状并对比离线生成与在线生成任务的差异，在此基础上明确本文的研究目标与主要贡献；

第二章给出在线实时手势生成任务的定义与系统设计原则，明确输入、输出模态的参数化表示，并介绍端到端实时数字人驱动系统的整体框架与模块划分；

第三章介绍本文多模态级联手势生成模型 FaceCapGes，包括头部姿态模态的引入方式、编码器结构设计以及身体姿态的层次化解码过程；

第四章介绍本文模型的训练方法，提出滑动窗口自回归训练与推理一致性策略，给出片段切割、窗口展开、监督损失与对抗训练目标的定义；

第五章介绍本文模型的主观与客观评估，并从定性分析、消融实验与性能测试等角度验证本文方法在生成质量与实时性方面的优势；

第六章总结全文工作并讨论未来研究方向。

第2章 多模态在线实时手势生成框架

本章围绕本文面向在线实时数字人驱动的手势生成任务，给出必要的概念定义、输入模态表示以及系统层面的总体设计框架。首先，我们明确本文研究所涉及的手势的定义与分类，并限定任务范围以满足严格低延迟场景的可行性。随后，本章介绍身体姿态与面部表情在计算机动画与深度学习中的参数化表示方式，为后续模型结构设计提供统一的输入输出形式。最后，本章给出端到端实时系统的整体架构与模块职责划分，说明 FaceCapGes 在实际部署系统中的位置与数据流组织方式。

2.1 手势的定义

在人类交流中，手势是常与语音同时出现的身体动作信号，它承担语义补充、情感表达与互动调节等多重功能。这类动作通常呈现出明确的表达意图，使其区别于姿态平衡，移动等纯功能性动作^[22]。因此，手势与语言并非两个相互独立的系统，而是源于共同的认知的表达机制。

在本研究的广义定义下，手势的运动形式不局限于手部运动，还可扩展至头部运动、躯干姿态等与表达相关的所有上半身动作。

2.1.1 Kendon 连续体

Kendon 连续体^[22-23]将与表达相关的手势行为，基于语言化程度做了以下分类：gesticulation（随语手势）→ language-like gestures（语言样手势）→ pantomime（拟态/哑剧式动作）→ emblems（约定俗成的象征手势）→ sign language（手语）。越靠左的类型通常更依赖当前的言语与语境、形式更即兴；越靠右则越接近离散的符号系统，规约化程度高，意义更稳定，可在缺少口语的情况下独立传达。

当前人机交互、虚拟人/数字人驱动等方向的手势生成，多数工作聚焦于 Kendon 连续体最左侧的随语手势，即说话过程中自然出现、与语音节奏与语义强关联的上肢动作。其含义往往依赖当前的口语与语境，脱离它们时通常难以传达清晰含义。

相对而言，连续体右侧的手势更接近文化中约定俗成的符号，可以脱离口语独立传达含义。例如，表达称赞的拍手动作，或表示“请保持安静”的嘴前竖起食指的动作，这些被视为象征手势，通常不在随语手势生成领域的研究对象中。

本文的研究范围据此限定在随语手势的学习与生成。

2.1.2 随语手势的分类

McNeill^[23]将随语手势进一步划分为四种基本类型：

1. **Iconic gestures (形象性手势)**: 以具象方式描绘事物的外形、空间路径或动作特征。例如，用手势勾勒一个物体的轮廓，或划出一道线表示移动的轨迹。此类手势与语言内容直接对应，表达具体语义。
2. **Metaphoric gestures (隐喻性手势)**: 表达抽象概念或思维结构的手势，比如用双手做出捧起一个物体放到一边的动作，表示“先把这个话题放一边”。这种手势并不描绘实体，而是以具象化的方式呈现抽象语义。
3. **Deictic gestures (指示性手势)**: 指向空间中的对象、人物或方向，常用于对话焦点的指明与注意引导。
4. **Beat gestures (节奏型手势)**: 与语音重音、韵律或节奏同步的节奏性动作，通常不承载具体语义，但可用于强调语音节奏，引起听众对说话内容的注意。

这四类手势构成了随语手势在语义与语篇功能上的主要维度，并在自然交流中常以复合形式出现。在手势生成任务中，研究通常将其视为不同的可生成目标：其中节奏型手势由于与语音韵律高度同步、对齐与建模相对容易，长期以来在数据驱动方法中更常被优先刻画；而形象性与隐喻性等语义相关手势则对文本的语义推理有更高要求，因而是更具挑战性的方向。

鉴于本文以低延迟在线驱动为目标，本文优先建模与韵律强耦合的节奏型手势；语义一致的形象性/隐喻性手势留作后续工作。

节奏型手势在随语手势中的重要性 尽管节奏型手势通常不承载具体语义信息，已有研究表明，其在交流效果与听众感知层面仍具有独立的价值。Baars 等的实验^[24]比较了无手势、仅使用节奏型手势、以及包含了形象性、隐喻性、指示性的意义性手势的三种演讲条件，结果显示，相较于完全不做手势，仅使用节奏型手势即可显著提升听众对说话者自然度的主观评价，并一定程度上提升了听众对演讲内容的记忆表现。而包含意义性手势的演讲条件在自然度与听众的记忆保留等指标上并未显著优于节奏型手势条件。这一发现表明，即便缺乏形象性或隐喻性的语义映射，节奏型手势仍能通过与语音韵律的同步，对交流过程产生积极影响。

从功能上看，节奏型手势主要服务于语篇结构与韵律组织，其作用并非传递附加语义，而是通过时间对齐、重音标记与注意力引导，增强语音信息的感知显著性与节奏感。在真实的人机交互与虚拟人系统中，这类手势常被作为一种低语义依赖、但高

度稳健的非语言表达形式加以采用。

鉴于本文面向低延迟、严格逐帧的在线驱动场景，系统在任一时间步均无法获取未来文本或完整语义结构，对语义一致性要求较高的形象性与隐喻性手势难以可靠生成。相比之下，节奏型手势主要依赖于当前及局部时间窗口内的语音韵律特征，更适合在实时条件下进行稳定建模与生成。因此，本文选择以节奏型手势作为主要研究对象，并将其视为一种在系统约束下具有明确交互价值的可行随语手势形式。

2.1.3 头部手势的分类

除手部动作外，头部动作同样是手势的重要组成部分。头部的点动与摆动在时间结构上常与手势及语音节奏保持同步^[9]，在语用功能上既能辅助语音韵律的组织，也能表达态度与指向信息。

在不同研究中，头部动作被从多个维度加以分析，其主要功能可归纳为以下几方面：

1. 韵律相关（**prosodic**）动作反映语音重音与句法节奏的对应关系^[25]；
2. 语义或态度相关（**semantic/attitudinal**）动作表达说话者的情绪倾向与交际意图^[22-23]；
3. 指向相关（**deictic**）动作通过转头或注视方向建立叙事空间的参照^[23]。

此外，研究表明头部动作的启动时间往往早于发声^[26]。具体而言，头部动作存在启动与加速过程，若其峰值需与重读音节的时间对齐，则动作必须提前起势。因此，头部动作可能对即将到来的语音韵律具有前瞻性。

这一特征揭示了头部动作与语音之间的时序关系，说明视觉模态中的运动信号有时可先于声学事件出现。本文研究也因此关注头部动作，并将其纳入输入模态。

2.2 身体姿态的参数化表示

手势作为身体运动的子集，其生成和识别依赖于身体姿态的连续建模。因此，在进一步讨论手势生成方法之前，在此明确身体姿态的参数化表示方式。

骨架结构的定义 在计算机动画与动作捕捉领域，身体姿态通常由骨架结构和关节旋转参数共同定义。骨架结构描述了人体各关节的拓扑关系及层级依赖；而每个姿态帧由一组关节旋转参数所确定，这些参数定义了相对于父节点的旋转变换。在不同的系统与任务中，骨架结构的具体形式往往有所差异，这种差异直接影响姿态数据的表

示与学习方式。

在不同的系统与任务中，骨架结构可以遵循各自的标准，因此，不同的数据集、3D 模型或神经网络往往基于自身定义的关节层级与命名体系进行训练与标注。例如，AMASS 数据集^[27]采用 SMPL 拓扑结构^[28]，BEAT 数据集^[5]使用简化上半身骨架。近年来的自动骨骼绑定与骨架归一化方法，通过学习或优化关节对应关系，实现了不同拓扑之间的姿态重定向（pose retargeting）^[29-30]，从而消除了模型依赖于特定骨架结构的限制。

旋转参数的选取 在确定骨架结构之后，具体的关节状态可通过多种旋转参数进行描述。常见的旋转参数表示方法包括：

1. 欧拉角（Euler Angles）：通过三个顺序旋转角表示姿态，直观且参数维度低（3 维），适合存储，但运算存在万向节锁（Gimbal Lock）问题，导致特定姿态下自由度丢失，且插值过程易产生非物理运动；
2. 四元数（Quaternion）：以四维单位向量 ($q_0^2 + q_1^2 + q_2^2 + q_3^2 = 1$) 表示旋转，避免欧拉角的奇异性的同时，实现平滑插值。然而在神经网络回归中通常需要显式处理单位长度约束，且同一旋转存在双覆盖（ q 与 $-q$ 等价），优化不稳定。
3. Axis-Angle：以旋转轴（三维单位向量）与旋转角度（标量）的乘积表示旋转，参数紧凑（3 维），是参数化人体模型 SMPL^[28]的主要姿态存储格式。但角度取值范围存在周期性 ($[0, 2\pi)$)，导致参数空间不连续，影响神经网络学习；
4. 旋转矩阵：以 3×3 的正交矩阵，表示坐标系或物体绕某个轴旋转一定角度后的新姿态。常用于计算机图形学、机器人学和物理仿真中，但正交约束 ($R^T R = I$ 且 $\det(R) = 1$) 提高神经网络的学习难度。
5. Rot6D^[10]：将 3×3 旋转矩阵的前两列展开为六维向量，通过 Gram-Schmidt 正交化过程自动保持列向量的正交性与单位长度，无需额外约束。该表示既继承了旋转矩阵的数值稳定性，又解决了其参数冗余和正交性约束难题，同时具备连续性优势。

在计算机图形与实时渲染中，通常采用四元数或旋转矩阵进行骨骼变换与插值，以保证数值稳定性和计算效率。然而，在深度学习任务中，这些表示方式的约束会对训练效率造成影响：四元数需满足单位长度约束；旋转矩阵则需满足元素的正交约束。近年来的研究表明^[10]，Rot6D 既规避了上述两种表示的约束难题，又保持了姿态空间的连续性与物理合理性，在动作生成与姿态预测任务中展现出更优的可学习性。

手势生成方法对旋转参数的选取经历了一个演变过程。早期工作采用的方法各

异（如欧拉角、Axis-Angle），多受限于特定框架或历史因素；而近期研究，出于对神经网络训练稳定性的追求，已显著趋向于采用 Rot6D 表示。这一趋势的实例对比见表 2.1。

表 2.1 不同旋转表示方式的空间连续性与使用示例

Table 2.1 Spatial Continuity and Usage Examples of Different Rotation Representations

表示方式	维度	连续性	典型应用场景	训练使用示例
欧拉角	3	存在万向节锁	姿态存储	CaMN ^[5]
四元数	4	连续，但存在单位长度约束	图形学	—
Axis-Angle	3	角度部分不连续	神经网络	DiffSHEG ^[14]
旋转矩阵	9	连续，但存在正交约束	图形学、机器人学	MambaGesture ^[31]
Rot6D	6	连续	神经网络	近期多数 ^[11-13]

因此，本文在姿态生成模型中统一使用 Rot6D 表示每个关节的旋转状态，以确保训练阶段的平滑收敛与推理阶段的稳定性。

2.3 面部表情的定义与参数化表示

面部表情是非语言交流的重要组成部分，与语音、手势共同传递情感和态度信息。与身体动作不同，面部表情主要由皮肤形变和局部运动构成，无法通过骨骼旋转直接建模，因此需要专门的参数化表示方式。

本文使用的表示方法为 BlendShape 权重模型。Facial Action Coding System (FACS)^[32] 提出复杂表情可分解为若干可组合的基本动作单元 (Action Units, AU)，为表情的参数化表示提供了理论基础。受此启发，BlendShape 模型将面部表情表示为一组可线性叠加的形变基，每个基形对应一个权重控制的局部表情变化。复杂表情由多个基形的组合生成，其思想与 FACS 的动作单元体系相呼应。

与基于骨骼的形变不同，BlendShape 直接在顶点层面定义几何偏移量，因此对网格拓扑结构高度依赖：每个基形的顶点偏移需与基础网格逐点对应。在具有相同网格结构的角色模型之间，BlendShape 集合可以直接复用；但若拓扑发生变化，偏移数据将无法一一对应，从而难以在不同模型间映射或重定向。这种拓扑依赖性限制了其跨模型的通用性。

尽管如此，BlendShape 在表现精细面部表情和软组织形变方面具有显著优势。其标准化权重接口、实时可驱动性与渲染兼容性，使其成为虚拟人和表情捕捉系统的主流表示形式，并被广泛应用于如 ARKit^[33] 等实时动画框架，以及主流渲染引擎中。

图 2.1 来自^[34]，展示 FACS 对闭眼睛的动作单元 AU45 的定义。

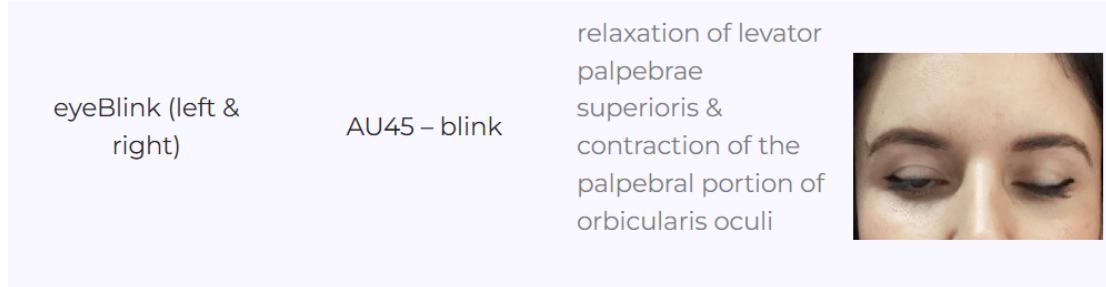


图 2.1 FACS 中闭眼动作单元的定义

Figure 2.1 Definition of Eye Closure Action Units in FACS

这在 BlendShape 中，一般对应两种基形：eyeBlinkLeft（闭左眼）与 eyeBlinkRight（闭右眼）。图 2.2 展示了 eyeBlinkLeft、eyeBlinkRight 的 ARKit BlendShape 基形在权重 $w \in [0, 1]$ 下的线性插值效果（从张眼到闭眼）。



图 2.2 BlendShape 在闭眼形变上的线性插值效果

Figure 2.2 Linear Interpolation of BlendShape Deformation for Eye Closure

2.4 系统设计思路

2.4.1 手势生成目标的选择

手势生成任务可根据其对语义上下文的依赖程度与时间结构的复杂度，大致划分为以语义理解为核心的“语义驱动型生成”，与以韵律对齐为核心的“韵律驱动型生成”。本文面向用户实时数字人驱动场景，在严格因果与低延迟约束下，模型无法

访问未来语音或完整文本语义，因此难以稳定生成对语义推理要求较高的形象性、隐喻性与指示性手势。基于该约束，本文将研究重点聚焦于与语音韵律高度同步的节奏型手势生成，并将其视为在实时交互条件下具有明确表达价值且可稳定建模的随语手势形式。

在此定位下，本文提出帧级多模态级联手势生成模型 FaceCapGes，输入为实时语音特征、面部 BlendShape 权重以及头部姿态旋转参数，并在不依赖未来帧信息的条件下逐帧输出上半身骨骼姿态。模型通过多模态信号融合弥补单一语音模态在实时场景下的信息不足：面部表情提供情绪与说话强度等辅助线索，头部姿态提供节奏前瞻与空间锚定信号，从而在保持可部署实时性的同时提升生成动作的自然度、同步性与方向一致性。

综合实时数字人驱动场景的交互需求与部署约束，本文系统设计遵循以下原则：

- (1) 严格因果性：模型仅使用当前与过去的多模态输入，不访问未来帧信息；
- (2) 低延迟：支持帧级在线推理，满足实时交互的时延要求；
- (3) 可实时采集模态：输入模态需能通过常见设备实时获取，包括语音、面部参数与头部姿态；
- (4) 可稳定建模：在上述约束下优先建模节奏型手势，并将其作为核心生成目标；

2.4.2 引入头部姿态的动机与贡献

从生成可行性的角度，现有研究普遍认为节奏型手势可在无语义理解的条件下由语音韵律直接驱动生成。多数语音驱动手势研究证实，仅凭语音的能量、时长与音高变化即可合成自然的节奏性上肢动作^[3,35-36]。这些研究所生成的动作在时间结构上与语音重音同步，体现了语音与手势共享的时间规划机制。

相比之下，iconic（形象性）、metaphoric（隐喻性）与 deictic（指向性）手势均依赖语义或指向关系，需要从上下文分析语义与语境，难以在严格实时的因果条件下生成。而节奏相关特征在音频中则具有更高的可预测性。^[4] 这表明，在缺乏未来语义与全局上下文的实时场景中，仅凭语音模态，模型只能稳定生成节奏层面的动作。

为突破这一限制，本文引入头部姿态模态作为补充输入信号。头部姿态能在实时因果条件下提供部分空间与时间线索：其转头与注视方向反映互动焦点，点头与抬头与语音重读共现，能够在不依赖未来语义信息的前提下，为手势生成提供弱先验约束。这种模态扩展为实时系统提供了理论上的可行性基础，使模型能够在语音之外获得关于节奏、方向与视角的附加信息。

2.4.2.1 头部姿态对手势预测的贡献

头部动作在自然语音中常呈现出一定的时间前瞻性^[26]: 其启动往往早于对应韵律词的发声, 这意味着视觉模态可能比声学信号更早反映语音节奏的变化趋势。这种时序特性为实时生成任务提供了潜在的预测窗口, 使系统能够在语音节奏变化尚未显现前, 就提前捕获相关的动态线索。因此, 头部姿态在实时生成中不仅提供同步参考, 也可能在时间上形成前驱信号, 为手势节奏的自然启动提供时序优势。

2.4.2.2 头部姿态对空间锚定与视角一致性的贡献

头部姿态模态为实时语音驱动的手势生成提供了关键的空间参照信号。其与语音韵律在时间组织上高度耦合。即使在无未来语义信息的条件下, 头部的转向与注视变化仍能反映说话者的注意焦点与叙述方向, 从而帮助模型在动作生成中保持空间的连贯性与方向一致性。这一机制使系统能够在时间与空间两个维度上同步对齐语音与动作, 让生成的手势在视觉上更具互动感与表达意图。

在 McNeill 的四类手势体系中, 头部姿态的引入主要强化了两类动作的生成:

- (1) 对 beat 手势而言, 它为语音重读和节奏段落提供显式的时间协同信号, 使手部与头部动作在韵律层面更加一致;
- (2) 对 iconic 手势而言, 它在具有路径与方向特征的动作中提供空间参考, 使模型能够在叙事空间中更稳定地确定动作的方位与轨迹方向。

通过这两方面的强化, 系统在保持实时性的同时获得了更自然的节奏衔接与空间表达。

与此同时, 本文明确头部姿态模态的作用边界: 其核心优势在于捕捉方向、焦点与时序节奏, 而非手型语义或复杂形态描摹等细粒度语义特征。换言之, 它主要改善手势的位置、方向与视角依附, 而非手势的形状描绘或语义内容。对于依赖抽象语义或外指参照的 metaphoric 与 deictic 手势, 仍需语言或上下文模态的补充。

总体而言, 头部姿态为实时生成提供了介于韵律与语义之间的关键中层约束。其时间上的前瞻性与空间上的指向性共同帮助模型在低延迟条件下保持自然、连贯且空间协调的动作表现, 从而在因果生成框架内有效拓展了语音驱动手势的可表达范围, 并为节奏主导型动作的实时生成提供了结构的支持。

基于上述任务范围与模态设计原则, 下一节将进一步给出本文系统的总体架构, 并说明各模块在实时生成流程中的功能定位。

2.5 在线实时手势生成系统整体框架

本节介绍整个系统的端到端驱动流程及模块职责划分。如图 2.3 所示，系统整体架构由五个层级组成：用户配置层、设备层、中间件层、手势生成模型层以及渲染与驱动层。各层之间通过多模态信号接口进行连接，实现从信号采集到虚拟人动作生成的端到端实时处理。

FaceCapGes 模型位于中间层，承担多模态输入到上半身姿态输出的核心推理任务，而输入采集与渲染模块分别负责信号获取与结果展示。

为实现基于语音、面部捕捉与头部姿态的实时数字人驱动系统，本文构建了完整的信号采集、动作生成与渲染展示的处理管线。FaceCapGes 模型作为该系统的核心计算模块，负责在实时约束下从多模态输入推理出当前帧的上半身骨骼姿态。

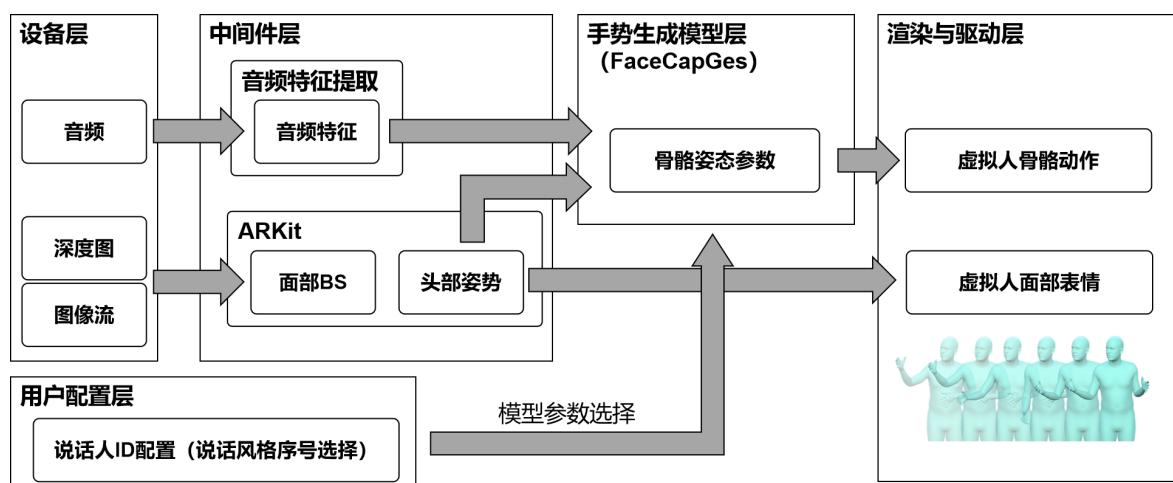


图 2.3 系统整体架构与数据流示意图
Figure 2.3 System Overview and Data Flow Diagram

2.5.1 信号采集与系统配置层

该层位于系统整体架构的输入端，用于从用户端设备实时获取多模态信号，并在系统初始化阶段完成运行参数的配置。整体结构可划分为设备层、中间件层与用户配置层三个部分，如图 2.3 所示。

设备层 设备层负责采集语音与视觉模态信号。语音信号由麦克风实时录制，采样率与帧移可根据运行设备性能调整；视觉信号由前置深度相机摄像头获取面部深度图与视频流，并作为 ARKit 面部追踪模块的输入。

中间件层 中间件层通过 Apple 提供的 ARKit 框架^[33]，将设备层的原始图像流与深度图转化为结构化特征。ARKit 输出两类主要数据：(1) **面部表情特征** ARKit 提供 52 维 BlendShape 系数向量，用于描述关键肌肉群的局部形变状态。该特征能够反映用户的表情、口型与情感变化，并以帧级形式同步输出。(2) **头部姿态特征** ARKit 在 ARFaceAnchor 中提供一个齐次变换矩阵 $\mathbf{T} \in \mathbb{R}^{4 \times 4}$ ，用于描述人脸锚点相对会话世界坐标系的位姿：

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix}, \quad \mathbf{R} \in \mathbb{R}^{3 \times 3}, \mathbf{t} \in \mathbb{R}^{3 \times 1}. \quad (2.1)$$

其中左上角的 \mathbf{R} 为旋转矩阵，右上角的 \mathbf{t} 为平移向量。本研究从矩阵中提取旋转部分 \mathbf{R} ，并将其转换为 Rot6D^[10]表示形式，以提升旋转空间的连续性与模型训练的稳定性。

同时，音频流在中间件层中被传入特征提取模块以生成时间序列特征。模型训练阶段使用 Librosa 库离线提取 Mel 频谱、短时能量与基频 F_0 等声学特征，以保证特征精度与一致性。系统运行阶段可由等价的实时特征提取模块（如 torchaudio 或 TensorFlow Audio）逐帧生成对应特征，以实现端到端的低延迟运行。

用户配置层 用户配置层负责系统初始化阶段的模型与参数设定。用户可在应用中选择说话风格，对应加载不同说话人 ID 配置下的模型权重。该配置仅在系统启动时生效，不参与实时推理过程。

本层提供的多模态信号经中间件处理后，以统一的数据接口传递至手势生成模型，实现语音、表情与头部姿态的实时融合输入。

2.5.2 手势生成模型层 (FaceCapGes)

FaceCapGes 模块作为系统的核心推理单元，接收来自信号采集与系统配置层的三类输入特征：语音特征、面部 BlendShape 系数以及头部姿态参数，并在不依赖未来帧的条件下，逐帧预测用户当前时刻的上半身骨骼姿态。

生成的骨骼姿态采用 Rot6D^[10]连续旋转表示形式，覆盖上半身 47 个关节的旋转参数。模型内部通过级联多模态编码结构提取时序相关特征，并利用单向 LSTM 解码器完成时间依赖建模，从而在保持实时性的同时，生成与语音节奏、表情变化及头部朝向高度一致的自然手势。

FaceCapGes 输出的姿态数据通过统一接口传递至渲染与驱动模块，与实时面部捕捉信号共同驱动虚拟角色的整体动作。由于模型仅依赖当前与历史帧输入，可与输入层以固定帧率并行运行，实现端到端的低延迟推理。

2.5.3 渲染驱动层

该模块位于系统输出端，负责将手势生成模型与面部捕捉结果共同转化为虚拟人的实时动作表现。系统将 FaceCapGes 模型输出的上半身骨骼姿态与 ARKit 实时检测的 52 维面部 BlendShape 系数传递至渲染引擎，由引擎内的模块解析并映射至目标虚拟人的骨骼与表情控制接口，从而实现多模态动作驱动。

渲染模块采用基于 GPU 的蒙皮计算与实时光照模型，以确保动画的平滑性和视觉一致性。最终，系统能够在实时流式输入条件下稳定运行，同步呈现语音、表情与身体动作，以自然流畅的数字人形象实现从多模态信号输入到可视化输出的完整驱动流程。

2.6 本章小结

本章围绕本文面向在线实时数字人驱动场景的多模态手势生成任务，给出了必要的概念定义、输入模态表示与系统层面的总体框架设计。首先，本章基于 Kendon 连续体与 McNeill 的随语手势分类体系，对本文研究对象的手势类型进行界定，并结合严格因果与低延迟约束，明确以节奏型手势作为可稳定建模的核心生成目标，为后续模型设计限定任务范围。随后，本章介绍了身体姿态与面部表情在动作生成任务中的参数化表示方式：姿态部分统一采用连续旋转表示 Rot6D，以保证训练与推理过程中的数值稳定性；面部表情部分采用 BlendShape 权重表示，并说明其与 ARKit 标准接口的对应关系。

在系统设计思路方面，本章进一步论证了引入头部姿态模态的必要性与贡献：头部姿态在实时因果条件下可提供节奏前瞻与空间锚定信号，从而弥补仅依赖语音模态时对方向一致性与互动焦点感知不足的问题，并在保持可部署实时性的前提下拓展手势生成的可表达范围。最后，本章给出了端到端在线实时系统的整体架构与模块职责划分，说明从多模态信号采集、特征转换到 FaceCapGes 推理与渲染驱动的完整数据流组织方式，为后续章节的模型结构设计与训练策略提供系统背景。

下一章将进一步详细介绍 FaceCapGes 的多模态级联模型架构，包括语音、面部与头部姿态特征的编码方式、融合策略以及面向上半身姿态的层次化解码过程。

第3章 融合头部姿态的多模态级联手势生成模型架构

CaMN^[5]提供了语音驱动手势生成的基础级联框架，为实时数字人动作建模奠定了有效的结构基础。然而，原始 CaMN 主要面向离线序列建模，其训练与推理过程默认可访问完整的时间上下文，难以直接满足在线实时场景下严格因果与逐帧生成的需求。为此，本文在保留级联解码思想的同时，对输入模态组织方式与时序建模进行了适配，使模型能够在仅依赖当前与历史信息的条件下进行生成。此外，为进一步增强空间方向感与节奏一致性，本文在语音与面部捕捉之外引入头部姿态模态作为额外输入，并设计对应编码器以提供补充表征。本章将介绍上述多模态编码、融合以及身体姿态解码的结构。

3.1 问题定义

在整体系统中，FaceCapGes 模块承担着从多模态输入信号到上半身骨骼姿态预测的核心任务。为了明确模型的输入输出结构与学习目标，本节对该问题进行形式化定义。

3.1.1 任务描述

目标是在实时条件下，根据用户当前时刻的语音、面部表情与头部姿态信息，预测其对应的上半身骨骼姿态。模型需能够逐帧生成与语音节奏、面部动态和头部转动方向相协调的自然手势动作，而不依赖未来的输入帧或整句语音信息。

形式上，可以将该任务定义为一个多模态时序映射函数：

$$\hat{\mathbf{v}}_t^B = f_{\theta}(\mathbf{v}_{t-N:t}^A, \mathbf{v}_{t-N:t}^F, \mathbf{v}_{t-N:t}^H), \quad (3.1)$$

其中 f_{θ} 表示由参数 θ 控制的生成模型， N 为历史窗口长度。各模态输入定义如下： \mathbf{v}_t^A 表示语音模态在时刻 t 的特征向量，由麦克风信号经特征提取模块得到； \mathbf{v}_t^F 表示面部模态的输入，为 ARKit 输出的 52 维标准化 BlendShape 系数； \mathbf{v}_t^H 表示头部模态的输入，为 ARKit 得出的头部旋转；而 $\hat{\mathbf{v}}_t^B$ 为生成模型在当前时刻预测的上半身骨骼姿态向量。模型仅利用当前及过去 N 帧的输入信息估计 $\hat{\mathbf{v}}_t^B$ ，从而满足严格的实时推理约束。

3.1.2 输入与输出模态

FaceCapGes 模型的输入由三种可同时实时获取的模态组成：语音特征、面部 BlendShape 权重及头部姿态参数；输出为当前帧的上半身骨骼旋转状态。各模态的符号与维度如表 3.1 所示。

表 3.1 输入输出模态符号与维度
Table 3.1 Notations and Dimensions of Input/Output Modalities

模态	符号	维度	描述
语音特征	v_t^A	\mathbb{R}^{1067}	由音频信号提取的时序特征（Mel 频谱、短时能量、基频等）
面部 BlendShape	v_t^F	\mathbb{R}^{52}	ARKit 输出的标准化表情权重向量
头部姿态	v_t^H	\mathbb{R}^6	头部关节旋转状态
骨骼姿态（输出）	\hat{v}_t^B	$\mathbb{R}^{6 \times 47}$	上半身 47 个关节的旋转状态

输入序列 $(v_{t-N:t}^A, v_{t-N:t}^F, v_{t-N:t}^H)$ 描述了用户在过去 N 帧内的语音与表情动态信息。模型通过学习其时序变化规律，逐帧生成对应的骨骼姿态输出 \hat{v}_t^B 。在推理阶段，模型仅访问至时刻 t 的输入序列，无法访问任何未来帧信息，保证了生成过程的因果性与实时性。

3.2 级联架构设计的继承

3.2.1 级联架构的原理与理论背景

现有语音驱动手势生成模型多采用多模态融合结构，其中以 CaMN^[5]为代表的级联架构在设计理念上具有代表性。其核心思想是将语音、面部表情与身体动作视为语义表达的不同层级：语音模态承担语义与节奏驱动作用，面部模态反映情感与意图，身体动作则是语言与情绪的外化呈现。CaMN 采用自上而下的处理顺序，即依次对语音、面部和动作模态进行建模，从而以层次化结构保持模态间的语义依存关系。

这种设计符合人类交流中“语言、表情、动作”一体化的认知规律^[22-23]。语音先规划语义与节奏，面部表情作为情绪强化信号随后产生，最终通过身体动作完成完整的非语言表达。模型中，语音编码器输出的时间嵌入被输入至面部编码器，再与面部特征融合后驱动动作解码器，从而保持语义一致性并增强表现力。

然而，CaMN 的原始设计面向离线整句生成任务，需要访问未来上下文以维持全局连贯性。在实时场景下，这种依赖将引入显著延迟并破坏因果性。FaceCapGes 在继承其层次思想的同时，对输入模式、训练方式与模态选择进行了系统性重构，以满足帧级实时约束。

3.2.2 说话人 ID 分支移除

如图 2.3 所示，用户配置层会设置说话人 ID 配置用于模型切换，但该模态在本模型中不属于网络输入。在基线模型 CaMN 中，输入模态包含显式的说话人 ID 向量，用于在同一模型内区分不同演讲者的风格差异。然而在实时交互场景下，该分支并非必要：用户身份通常固定，且说话风格的变化频率远低于帧级推理速度。因此，FaceCapGes 移除了 ID 输入分支，采用针对每个说话人独立训练模型参数的方法。实验表明，该方式能在保持收敛稳定的同时提升动作的自然性与节奏一致性。从系统使用角度看，不同模型可视为说话风格配置，用户仅在需要时切换对应参数，该操作发生频率低，不会影响实时推理性能。

3.2.3 输入模态的解码

语音特征通过时间卷积网络（Temporal Convolutional Network, TCN）和多层次感知机（Multilayer Perceptron, MLP）编码，以捕捉短时节奏模式；面部模态采用相似结构，并在中间层融合语音嵌入，从而增强语音与表情之间的语义关联。两者都延续 CaMN^[5]的结构设计。

具体而言，语音编码器 E_A ，采用 12 层 TCN 建模局部时间依赖，并通过跳跃连接（skip connection）增强深层特征的传递稳定性；在第 12 层提取的语音时序特征后接入两层 MLP，用于进一步特征精炼与维度压缩，最终输出语音潜在表示 $z_t^A \in \mathbb{R}^{128}$ 。面部编码器 E_F 采用较浅的 8 层 TCN 结构以捕捉面部表情的短时动态变化，并在第 8 层将语音嵌入与面部特征进行通道维拼接融合，以增强面部表情变化与语音节奏之间的对应关系；编码器末端同样使用两层 MLP 进行特征映射与输出压缩，最终得到面部潜在表示 $z_t^F \in \mathbb{R}^{32}$ 。

语音编码器 E_A 与面部编码器 E_F 的输出定义为：

$$z_t^A = E_A(\mathbf{v}_{t-N:t}^A), \quad z_t^F = E_F(\mathbf{v}_{t-N:t}^F; z_t^A). \quad (3.2)$$

这两个编码器负责提取低层次语音节奏与表情动态信息，为后续模态融合提供稳定上下文表征。

此外，系统在此基础上引入头部姿态模态 \mathbf{v}_t^H ，用于补充空间方向与节奏信号。其编码器 E_H 将头部旋转向量映射为紧凑潜在表征：

$$z_t^H = E_H(\mathbf{v}_{t-N:t}^H), \quad (3.3)$$

编码器结构将在第 3.3 节详细说明。

3.2.4 身体姿态的解码

在输入模态经过编码与融合后，模型需将多模态特征映射至对应的身体姿态空间。为实现层次化的动作生成与结构协调，本文将上半身的输出区域划分为两个互补分支：躯干（Torso, T）与上肢（Upper limbs, U）。躯干部分包含脊椎的三个主要控制关节，用于确定身体的姿态基准与运动节奏；上肢部分包含双臂及手部关节，负责生成与语音节奏及情绪表达相呼应的细节动作。最终的上半身姿态表示为两者的组合：

$$\mathbf{v}_t^B = \mathbf{v}_t^T \otimes \mathbf{v}_t^U, \quad (3.4)$$

其中 \otimes 表示通道维度拼接操作。

该分层设计继承了 CaMN 的层次预测思路：模型首先生成相对稳定的躯干姿态以确定整体方向，再以此为条件预测上肢动作，从而在实时生成中保持整体协调性与自然度。

融合输入的序列化表示 在时刻 t ，来自语音、面部与头部编码器的特征 \mathbf{z}_t^A 、 \mathbf{z}_t^F 、 \mathbf{z}_t^H 与历史动作上下文共同构成 LSTM 解码器的输入。为保持严格因果性，我们显式提供最近 N 帧的历史上半身姿态 $\mathbf{v}_{t-N:t-1}^B$ 作为条件信息，并对当前时刻的动作输入使用占位符进行对齐。具体而言，定义历史动作对齐序列 \mathbf{s}_τ 为：

$$\mathbf{s}_\tau = \begin{cases} \mathbf{v}_\tau^B, & \tau \leq t-1, \\ \mathbf{0}, & \tau = t, \end{cases} \quad \tau \in \{t-N, \dots, t\}, \quad (3.5)$$

则窗口内每一帧的融合输入向量可写为：

$$\mathbf{z}_\tau^{fuse} = \mathbf{z}_\tau^A \otimes \mathbf{z}_\tau^F \otimes \mathbf{z}_\tau^H \otimes \mathbf{s}_\tau, \quad \tau \in \{t-N, \dots, t\}, \quad (3.6)$$

并将其按时间维堆叠得到长度为 $N+1$ 的因果上下文输入序列：

$$\mathbf{Z}_t^{fuse} = (\mathbf{z}_{t-N}^{fuse}, \dots, \mathbf{z}_t^{fuse}). \quad (3.7)$$

躯干与上肢的级联解码 将输入序列 \mathbf{Z}_t^{fuse} 分别送入躯干与上肢两路单向 LSTM 解码器，得到窗口内的输出序列：

$$\mathbf{O}_t^T = \text{LSTM}_T(\mathbf{Z}_t^{fuse}), \quad \mathbf{O}_t^U = \text{LSTM}_U(\mathbf{Z}_t^{fuse}), \quad (3.8)$$

其中 $\mathbf{O}_t^T = (\mathbf{o}_{t-N}^T, \dots, \mathbf{o}_t^T)$ ， $\mathbf{O}_t^U = (\mathbf{o}_{t-N}^U, \dots, \mathbf{o}_t^U)$ 。由于本文在时刻 t 的目标是预测当前帧动作，我们仅取序列末端输出作为当前帧的潜在表征：

$$\mathbf{z}_t^T = \mathbf{o}_t^T, \quad \mathbf{z}_t^U = \mathbf{o}_t^U. \quad (3.9)$$

最后，通过两路独立的 MLP 模块将潜在表征还原为旋转参数：

$$\hat{\mathbf{v}}_t^T = \text{MLP}_T(\mathbf{z}_t^T), \quad \hat{\mathbf{v}}_t^U = \text{MLP}_U(\mathbf{z}_t^U), \quad (3.10)$$

并拼接得到当前帧的完整上半身动作预测：

$$\hat{\mathbf{v}}_t^B = \hat{\mathbf{v}}_t^T \otimes \hat{\mathbf{v}}_t^U. \quad (3.11)$$

上述内容描述了模型在一个因果上下文输入序列内完成对当前帧动作的解码过程。

需要说明的是，LSTM 在实现中维护隐藏状态与记忆单元状态以编码时间依赖，其形式化表达见第4.1节。

3.3 头部姿态模态的引入

3.3.1 输入特征的定义

本文仅使用头部旋转信息作为输入特征，不引入头部位置或位移。

这是因为，头部位置相对旋转易受到身体姿态的变化。以 BEAT^[5]为例的演讲数据集中，演讲者多为站姿录制，单次演讲时长偏长（约 1 分钟），有轻微的重心移动等站姿调整；而目标应用场景下的用户交互姿态有可能包含坐姿，此时的身体活动方式与站姿存在差异，例如不存在长时间的站姿带来的重心移动。由此可见，直接建模头部位置或位移可能引入额外的场景依赖性，从而削弱跨场景泛化能力。因此，本文仅采用头部旋转作为输入。

3.3.2 特征获取方法

如第 2.5 节所述，头部姿态在推理阶段可从面部捕捉工具（如 ARKit^[33]）实时提取。但在训练阶段中，演讲数据集通常不单独提供头部姿态信号。此时，我们可以在预处理中利用骨架层级关系，沿骨架链路从根节点到头部关节的相对旋转进行旋转姿态的叠加，从而得到头部在全局坐标系下的绝对旋转表示。

3.3.3 级联架构中的位置

在级联架构设计中，我们考察了头部姿态特征与其他模态的多种组合方式。具体而言，分别尝试了：(1) 将头部姿态特征在编码阶段与语音或面部特征进行早期融合；(2) 在解码阶段以前两者的嵌入结果为条件，预测头部姿态特征作为辅助信号。实验结果显示，这两种交互方式均未带来显著性能提升，观察到了训练收敛速度的下降。

本文推测，头部动作虽然与语音韵律在时间上存在同步性，但头部动作与语音或表情之间的驱动关系不强，因为头部动作可以包含演讲人自然面向不同方向的听众等，与语音韵律或情感难以找到关联性的信息。

基于此观察，本文在最终架构中采用了弱耦合的设计：头部姿态特征在语音与面部特征编码完成后，以独立通道的形式拼接至多模态隐向量 \mathbf{z}_t^{fuse} ，同时头部姿态特征不参与其它输入模态的编码过程。我们在此配置下得到了更快的训练收敛。

编码器结构 图 3.1 所示为头部姿态编码器结构。该编码器由两层前馈网络组成，输入为 Rot6D^[10]表示的 6 维向量：

$$\mathbf{z}_t^H = E_H(\mathbf{v}_{t-N:t}^H), \quad (3.12)$$

其中 E_H 的具体形式为：

$$\mathbf{h}_1 = \text{ReLU}(\mathbf{W}_1 \mathbf{v}_t^H + \mathbf{b}_1), \quad (3.13)$$

$$\mathbf{z}_t^H = \mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2, \quad (3.14)$$

网络维度设置为：输入 6，中间层 36，输出 12。在特征层面，其输出与语音、面部嵌入拼接后输入解码器，形成从语义到反应的多层次信号流。



图 3.1 头部姿态编码器结构示意图
Figure 3.1 Architecture of the Head Pose Encoder

3.4 模型整体结构

图 3.2 展示了 FaceCapGes 从音频、面部、头部编码器分别提取模态特征后拼接，输入至 LSTM 解码器生成躯干与手部动作的过程。其中，训练阶段历史姿态序列比目标长度少一帧，需进行零填充进行对齐，如式 (3.6) 所示。

3.5 本章小结

本章围绕 FaceCapGes 的模型结构设计，系统介绍了面向严格实时交互场景的多模态级联架构与关键输入输出模态设定。首先，本文在任务设定上明确了逐帧在线生成与严格因果约束：模型在任意时刻仅利用当前及过去帧的语音、面部与头部信息进

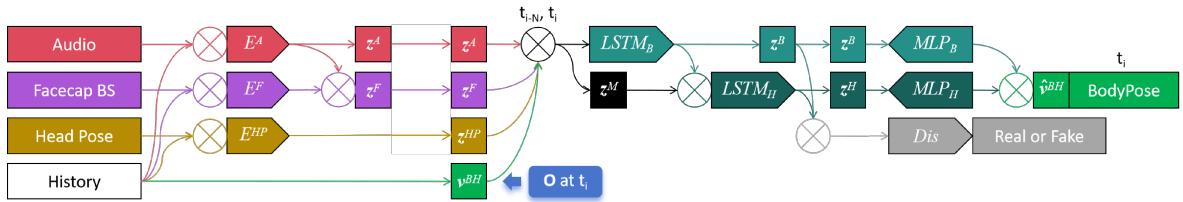


图 3.2 FaceCapGes 模型整体结构
Figure 3.2 Overall Architecture of FaceCapGes

行动作预测，从而区别于依赖未来上下文的一次性离线生成方法，并为后续结构设计提供了统一前提。

在模型结构方面，本文继承 CaMN 的级联框架，将语音、面部表情与身体动作视为具有层次关联的表达信号，并进一步系统性地引入头部姿态作为独立输入模态，以补充空间方向与节奏信息。同时，本文仅使用头部旋转信息进行建模，从而提升跨场景泛化能力，并给出了分别适用于训练与推理阶段的特征获取方式。

在解码端，本章说明了面向上半身骨骼姿态的层次化输出设计：将动作划分为躯干与上肢两路分支，并通过级联解码保持整体协调性与自然度。上述模型结构为后续章节提出的滑动窗口自回归展开策略与片段级优化目标提供了清晰的结构基础。下一章将进一步介绍模型在训练与推理阶段的统一展开过程，并证明该策略在严格因果约束下能够实现稳定的实时动作生成。

第4章 基于滑动窗口的实时手势生成自回归训练

本章讨论在严格因果约束下，如何将第3章定义的窗口内单步预测模型展开为逐帧自回归生成过程，并构建与之对应的训练目标。由于实时场景无法观测未来的多模态输入，模型必须在每个时间步仅依赖历史信息预测下一帧动作，从而保证在线推理的一致性与可部署性。为此，本文提出片段切割与滑动窗口展开的流程：通过前置动作帧缓冲历史上下文，并在生成阶段逐帧预测与写回动作缓存，形成连续的片段级输出序列。在此基础上，本章进一步给出片段级监督损失与对抗训练目标的定义。

4.1 单步因果 LSTM 预测器

为实现严格因果的实时手势生成，本章首先建立一个用于“下一帧动作预测”的基本单元：单步因果 LSTM 预测器。该预测器在每个时间步利用当前可用的多模态输入特征，并结合历史动作上下文与循环状态，对下一帧动作进行估计。在后续章节中，我们将该单步预测器以滑动窗口方式展开，从而形成对一个动作片段的自回归生成过程，并据此定义片段级别的监督损失与对抗训练目标。

4.1.1 LSTM 基本概念与状态传递机制

手势动作序列具有显著的时间依赖性：当前姿态不仅受当前输入模态（如语音、面部表情、头部姿态等）影响，也与过去的动作状态密切相关。循环神经网络（Recurrent Neural Network, RNN）通过引入随时间递推的隐状态来建模序列依赖，而 LSTM 进一步通过门控机制缓解长序列训练中的梯度消失问题，从而更适合用于动作序列建模。

在标准的 LSTM 结构中，网络在每个时间步接收当前输入特征 \mathbf{x}_t ，并维护两类递推状态：隐藏状态 \mathbf{h}_t 与记忆单元状态 \mathbf{c}_t 。其中， \mathbf{h}_t 可视为与当前输出相关的短期表示，而 \mathbf{c}_t 作为更稳定的记忆轨道，用于在更长时间尺度上保留信息。其递推过程可形式化表示为：

$$(\mathbf{h}_t, \mathbf{c}_t) = \text{LSTM}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}). \quad (4.1)$$

根据隐状态传播方向的不同，LSTM 主要分为单向 LSTM（Unidirectional LSTM）与双向 LSTM（Bidirectional LSTM）。

双向 LSTM 通过同时构建前向与反向递推路径，在输出时刻 t 的表示时会融合来自未来时间步的反向信息。该结构在离线分析或全序列可见的任务中能够更充分

利用上下文，从而提升预测性能；然而在实时生成任务中，未来输入在时刻 t 尚不可用，反向路径的依赖无法满足。

相比之下，单向 LSTM 沿时间正向递推，其状态更新仅依赖于过去与当前输入，从结构上满足因果约束。对于下一帧动作生成问题，UniLSTM 能够自然地被解释为一个逐步更新的预测器：在每个时间步基于当前可观测输入与历史状态输出下一帧动作估计，并将生成结果反馈至下一步，从而形成自回归（autoregressive）的实时生成流程。此外，单向递推的状态传递机制也使得模型能够在有限的输入上下文之外保留更长程的动态信息，为后续滑动窗口展开提供了必要的记忆能力。

基于以上原因，本文采用单向 LSTM 作为动作生成的时序建模骨干网络，以满足严格因果的实时生成需求。在后续的滑动窗口展开策略中（见第4.3节），该单步预测器将被逐帧迭代调用，从而在片段级别完成动作序列的自回归生成，并用于定义监督损失与对抗训练目标。

4.1.2 双时间尺度记忆结构

尽管单向 LSTM 通过递推状态 $(\mathbf{h}_t, \mathbf{c}_t)$ 在理论上具备建模长程依赖的能力，但在自回归动作生成任务中，若仅依赖循环状态作为唯一的历史信息通道，将迫使该状态同时承担短期细节与长期规律的表征。由于循环状态本质上是对历史信息的压缩表示，其容量有限且受 LSTM 的门控遗忘机制影响，短期运动细节与局部连续性约束可能难以稳定保留，从而使生成过程更易出现抖动、漂移或长期一致性下降等现象。因此，本文在跨步循环状态之外显式提供固定长度的历史动作上下文，作为短期约束信号，以在保持长期记忆能力的同时增强逐帧预测的稳定性与可控性。

为提升逐帧生成的稳定性，本文在循环状态之外显式引入固定长度的历史动作上下文作为短期条件信息。具体而言，设 N 为历史动作上下文长度，在预测时刻 t 的动作时，我们显式提供最近 N 帧的动作历史序列 $\mathbf{v}_{t-N:t-1}^B$ 作为可观测输入，使模型在每一步预测中均能获得短期运动细节与局部连续性约束。为保持因果性，该历史动作序列不包含当前待预测时刻 t 的真实动作，而采用占位符进行对齐（例如以零向量或掩码符号表示），从而形成长度为 $N+1$ 的因果上下文序列。

在该设计下，模型的历史信息将通过两条互补路径传递：一方面，循环状态 (\mathbf{h}, \mathbf{c}) 作为长期记忆通道，用于编码超过 N 帧范围的更长程动态趋势、说话风格与运动节奏；另一方面，显式历史动作上下文作为短期精确条件，在每一步预测中直接提供最近 N 帧的局部运动信息，从而缓解状态漂移带来的不确定性，并提升自回归生成的鲁棒性。本文采用“短期显式上下文 + 长期隐式状态”的双时间尺度记忆结构，以在

严格因果约束下平衡表达能力与生成稳定性。

4.1.3 窗口内预测过程的形式化表达

本节将第3章定义的窗口内解码结构形式化为可递推的单步因果预测器，以明确在线生成时循环状态的跨步传递方式。在时刻 t ，模型以长度为 $N+1$ 的因果上下文输入序列 \mathbf{Z}_t^{fuse} 作为输入（其定义见式(3.7)），并分别维护躯干与上肢的 LSTM 隐藏状态与记忆单元状态： $(\mathbf{h}_{t-1}^T, \mathbf{c}_{t-1}^T)$ 与 $(\mathbf{h}_{t-1}^U, \mathbf{c}_{t-1}^U)$ 。

在线运行时，我们将跨步传递的状态作为 LSTM 解码器的初始状态，从而得到窗口内输出序列并更新状态：

$$\mathbf{O}_t^T, (\mathbf{h}_t^T, \mathbf{c}_t^T) = \text{LSTM}_T(\mathbf{Z}_t^{fuse}, \mathbf{h}_{t-1}^T, \mathbf{c}_{t-1}^T), \quad (4.2)$$

$$\mathbf{O}_t^U, (\mathbf{h}_t^U, \mathbf{c}_t^U) = \text{LSTM}_U(\mathbf{Z}_t^{fuse}, \mathbf{h}_{t-1}^U, \mathbf{c}_{t-1}^U). \quad (4.3)$$

由于本文在时刻 t 的目标是预测当前帧动作，我们使用窗口末端输出作为当前帧潜在表征，并通过 MLP 解码得到 $\hat{\mathbf{v}}_t^T$ 与 $\hat{\mathbf{v}}_t^U$ ，最终拼接为 $\hat{\mathbf{v}}_t^B$ 。该过程已在第3.2.4节给出，此处不再赘述。

因此，可以将单步因果预测器抽象为如下递推形式：

$$\hat{\mathbf{v}}_t^B, (\mathbf{h}_t, \mathbf{c}_t) = f_\theta(\mathbf{Z}_t^{fuse}, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}), \quad (4.4)$$

其中 $(\mathbf{h}_t, \mathbf{c}_t)$ 表示躯干、上肢 LSTM 状态的集合。在第4.3节中，我们将进一步描述该单步预测器如何随时间滑动展开，从而在片段尺度上生成长度为 M 的动作序列，并用于定义监督损失与对抗训练目标。

4.2 训练片段切割

本工作沿用 CaMN 采用的训练样本构造方式^[5]，将长序列动作数据切割为固定长度的短片段作为训练样本。

4.2.1 固定长度片段定义

设原始动作序列为 $\{\mathbf{v}_t^B\}_{t=1}^T$ ，以及对应的多模态输入特征序列 $\{\mathbf{z}_t^A, \mathbf{z}_t^F, \mathbf{z}_t^H\}_{t=1}^T$ 。我们从长序列中截取长度为 L 的连续片段作为训练样本，其中片段长度由历史上下文长度 N 与片段内生成步数 M 共同决定：

$$L = N + M. \quad (4.5)$$

在本文设定中，沿用 CaMN 的片段长度配置取 $L = 34$ 帧，同时选取历史上下文长度 $N = 16$ 帧，因此对应的片段内自回归生成步数为 $M = 18$ 帧。

对于第 k 个训练片段，其时间范围为 $[s_k, s_k + L - 1]$ ，片段内的动作与多模态输入分别表示为：

$$\mathbf{V}_k^B = (\mathbf{v}_{s_k}^B, \mathbf{v}_{s_k+1}^B, \dots, \mathbf{v}_{s_k+L-1}^B), \quad (4.6)$$

$$\mathbf{Z}_k = \left(\{\mathbf{z}_{s_k:t}^A\}_{t=s_k}^{s_k+L-1}, \{\mathbf{z}_{s_k:t}^F\}_{t=s_k}^{s_k+L-1}, \{\mathbf{z}_{s_k:t}^H\}_{t=s_k}^{s_k+L-1} \right). \quad (4.7)$$

其中 \mathbf{V}_k^B 将进一步划分为两部分：前 N 帧作为片段的历史上下文（亦即前置动作帧，用于缓冲与提供因果条件），后 M 帧为片段内需要逐帧生成并参与训练目标计算的部分。该划分将于第4.3节中用于描述滑动窗口展开过程与片段级损失计算策略。

4.2.2 重叠切割与样本覆盖

为提高训练样本覆盖率并增强模型对不同对齐位置的鲁棒性，我们采用带重叠的滑动切割方式从长序列中提取片段。具体而言，片段起始位置 s_k 按固定步长 Δ 滑动：

$$s_k = 1 + (k - 1)\Delta, \quad (4.8)$$

从而得到一组相互重叠的训练片段 $\{\mathbf{V}_k^B, \mathbf{Z}_k\}$ 。在本文实现中，沿用 CaMN 的设置取 $\Delta = 10$ 帧，以在样本数量与数据冗余之间取得平衡。

4.3 片段内部的滑动窗口展开策略

上一节定义了训练样本以固定长度片段组织：每个片段长度为 $L = N + M$ ，其中前 N 帧为历史上下文，后 M 帧为需要逐帧生成的目标区间。本节进一步描述片段内部的滑动窗口展开策略：该策略在训练与推理阶段保持一致，以单步因果预测器（第4.1节）为基本计算单元，逐帧生成长度为 M 的动作序列并拼接得到片段级输出。在本文设定中， $L = 34$ ， $N = 16$ ，因此 $M = 18$ ，窗口长度为 $N+1 = 17$ 。

4.3.1 前置动作帧与预热阶段

在严格因果的实时生成中，模型在时刻 t 预测动作时不能访问未来输入，并且其输出需要在片段之间保持连续。为此，我们将每个片段的前 N 帧动作视为前置动作帧，其作用是为后续生成提供短期运动上下文，并避免片段边界处出现断裂。

在片段开始时刻 s_k ，我们首先执行预热阶段：将片段的前 N 帧真实动作 $\{\mathbf{v}_{s_k}, \dots, \mathbf{v}_{s_k+N-1}\}$ 写入历史动作缓存 \mathcal{H} ，并同步读取对应的多模态输入特征。本文在预热阶段不执行

任何前向计算，生成阶段的初始 LSTM 状态由预设初始化值给出。本文采用零初始化作为 $(\mathbf{h}_{t-1}, \mathbf{c}_{t-1})$ 的初值。

4.3.2 滑动窗口展开与逐帧自回归生成

在预热阶段结束后，我们进入生成阶段：模型在片段内部进行 M 步滑动窗口展开，每一步生成 1 帧动作，并将预测结果写回历史动作缓存以供下一步使用。设片段内生成阶段的第 m 步对应全局时间 $t = s_k + N - 1 + m$ （其中 $m = 1, 2, \dots, M$ ），历史动作缓存 \mathcal{H}_{t-1} 包含最近 N 帧可用动作序列：

$$\mathcal{H}_{t-1} = (\tilde{\mathbf{v}}_{t-N}^B, \dots, \tilde{\mathbf{v}}_{t-1}^B), \quad (4.9)$$

其中 $\tilde{\mathbf{v}}$ 表示当前可用的动作帧：在生成开始时它包含真实前置动作帧，在生成推进过程中则逐渐由模型预测结果覆盖。

在时间步 t ，我们构造长度为 $N+1$ 的因果上下文窗口，组合当前可观测的多模态输入特征与历史动作缓存，形成融合输入序列 \mathbf{Z}_t^{fuse} （其定义见式(3.7)）。其中历史动作序列部分采用 $(\tilde{\mathbf{v}}_{t-N}^B, \dots, \tilde{\mathbf{v}}_{t-1}^B, \mathbf{0})$ 进行对齐，末帧使用占位符以保证严格因果性。

随后，单步因果预测器以 \mathbf{Z}_t^{fuse} 与跨步 LSTM 状态作为输入，输出当前帧动作预测 $\hat{\mathbf{v}}_t^B$ 并更新状态（窗口内形式化表达见第4.1.3节）。生成结果 $\hat{\mathbf{v}}_t^B$ 会被写回历史动作缓存，从而更新 \mathcal{H}_t 并用于下一时间步预测：

$$\mathcal{H}_t = (\tilde{\mathbf{v}}_{t-N+1}^B, \dots, \tilde{\mathbf{v}}_{t-1}^B, \hat{\mathbf{v}}_t^B). \quad (4.10)$$

通过上述逐帧自回归展开，我们得到片段内生成区间的预测序列：

$$\hat{\mathbf{V}}_k^{gen} = (\hat{\mathbf{v}}_{s_k+N}^B, \hat{\mathbf{v}}_{s_k+N+1}^B, \dots, \hat{\mathbf{v}}_{s_k+L-1}^B), \quad (4.11)$$

其长度为 M 。

4.3.3 拼接生成序列与片段级输出

由于滑动窗口展开在每一步仅生成 1 帧动作，片段级输出通过对 M 步预测结果进行时间维拼接获得。片段级生成结果 $\hat{\mathbf{V}}_k^{gen}$ 与真实动作片段 \mathbf{V}_k^B 在时间上对齐，其中前 N 帧为可观测上下文，后 M 帧为模型生成输出。因此，后续训练目标的定义将以 $\hat{\mathbf{V}}_k^{gen}$ 为核心对象，并与真实片段中对应的生成区间进行比较。

4.4 监督损失

在训练阶段，给定来自多模态语音动作数据集的配对样本序列：

$$(z_t^A, z_t^F, z_t^H, v_t^B), \quad (4.12)$$

模型的学习目标是在严格因果约束下生成与输入相匹配的上半身动作序列。与传统离线序列预测不同，本文采用第4.3节所述的滑动窗口展开策略：模型在每个片段内部自回归地逐帧生成 M 帧动作，并将其拼接为片段级生成序列。

4.4.1 片段级生成序列与损失计算范围

对于第 k 个训练片段，其长度为 $L = N + M$ （第4.2节），其中前 N 帧为前置动作帧，后 M 帧为模型需要生成的目标区间。根据滑动窗口展开过程（式(4.11)），模型得到片段生成区间的预测序列：

$$\hat{\mathbf{g}}_k = (\hat{\mathbf{v}}_{s_k+N}^B, \hat{\mathbf{v}}_{s_k+N+1}^B, \dots, \hat{\mathbf{v}}_{s_k+L-1}^B) \in \mathbb{R}^{M \times d}, \quad (4.13)$$

其中 d 表示动作表示的维度（本文中使用 Rot6D^[10]，即 $d = 6$ ）。对应的真实动作序列为：

$$\mathbf{g}_k = (v_{s_k+N}^B, v_{s_k+N+1}^B, \dots, v_{s_k+L-1}^B) \in \mathbb{R}^{M \times d}. \quad (4.14)$$

需要强调的是，片段前 N 帧前置动作仅用于提供因果历史上下文与片段平滑过渡，其本身并非生成目标。因此，与 CaMN^[5] 不同的是，本文的监督损失仅在片段生成区间 $\{s_k + N, \dots, s_k + L - 1\}$ 上计算，不对前置动作帧计算任何损失项。这是因为前置动作帧作为已知条件用于初始化历史缓存，而模型输出仅对应后续的生成区间。

4.4.2 总体优化目标

综合考虑空间重构精度、时序平滑性以及动作分布一致性，本文的总体优化目标定义为：

$$\mathcal{L}_{total} = \lambda_r \mathcal{L}_{rec} + \lambda_v \mathcal{L}_{vel} + \lambda_a \mathcal{L}_{acc} + \lambda_{adv} \mathcal{L}_{adv}, \quad (4.15)$$

其中 \mathcal{L}_{rec} 衡量片段生成区间的姿态重构误差， \mathcal{L}_{vel} 与 \mathcal{L}_{acc} 分别约束速度与加速度的连续性， \mathcal{L}_{adv} 表示对抗训练损失，将在第4.5节中进一步介绍。

4.4.3 姿态重构与时序平滑损失

为同时保证空间重构精度与时间连续性，我们采用基于 Huber 误差的重构损失形式，并分别作用于姿态、速度与加速度信号。给定任意预测序列 $\hat{\mathbf{x}}$ 及其对应的真实

序列 \mathbf{x} , 基础误差项定义为:

$$\mathcal{L}_{Huber}(\mathbf{x}, \hat{\mathbf{x}}) = \beta \cdot \text{SmoothL1}\left(\frac{\mathbf{x}}{\beta}, \frac{\hat{\mathbf{x}}}{\beta}\right), \quad (4.16)$$

其中 SmoothL1(\cdot) 表示平滑 L1 误差, β 为平滑系数, 本文中设为 0.1。

在此基础上, 片段生成区间的姿态、速度与加速度损失分别定义为:

$$\mathcal{L}_{rec} = \mathcal{L}_{Huber}(\mathbf{g}_k, \hat{\mathbf{g}}_k), \quad (4.17)$$

$$\mathcal{L}_{vel} = \mathcal{L}_{Huber}(\mathbf{g}'_k, \hat{\mathbf{g}}'_k), \quad (4.18)$$

$$\mathcal{L}_{acc} = \mathcal{L}_{Huber}(\mathbf{g}''_k, \hat{\mathbf{g}}''_k), \quad (4.19)$$

其中一阶与二阶时间差分 \mathbf{g}'_k 、 \mathbf{g}''_k 定义为:

$$\mathbf{g}'_{k,t} = \mathbf{g}_{k,t} - \mathbf{g}_{k,t-1}, \quad \mathbf{g}''_{k,t} = \mathbf{g}'_{k,t} - \mathbf{g}'_{k,t-1}, \quad (4.20)$$

预测序列 $\hat{\mathbf{g}}'_k$ 、 $\hat{\mathbf{g}}''_k$ 同理定义。

上述多尺度重构约束在自回归预测过程中能够有效缓解高频抖动与速度漂移问题, 在保证运动学精度的同时提升生成序列的时间稳定性。

4.4.4 损失权重设置

各损失项的权重系数在实验中设定为 $\lambda_r = 5 \times 10^2$, $\lambda_v = 10^3$, $\lambda_a = 10^3$, $\lambda_{adv} = 10^{-1}$ 。

4.5 对抗训练

尽管第4.4节的监督损失能够约束生成序列在逐帧空间误差与局部平滑性上的一致性, 但仅依赖点对点重构目标往往难以完全刻画真实动作序列的整体动力学分布。为进一步提升生成动作的自然度与分布一致性, 本文引入片段级判别器, 在片段尺度上约束生成序列与真实序列的统计特性一致。

4.5.1 基于拼接片段的片段级判别

与窗口内部的中间输出不同, 本文的判别器直接以片段生成区间的拼接序列为输入。对于第 k 个片段, 生成器输出的预测序列 $\hat{\mathbf{g}}_k$ 及其对应的真实序列 \mathbf{g}_k 定义见式(4.13)与式(4.14), 二者均为长度 M 的序列。判别器 $Dis(\cdot)$ 接收一段动作序列, 并输出其来自真实数据分布的概率:

$$Dis(\mathbf{g}) \in (0, 1). \quad (4.21)$$

通过片段级输入, 判别器能够从整体动力学角度判断生成动作的真实感, 从而在节奏、能量变化与运动统计特性等层面提供补充监督信号。

前置动作帧的掩码策略 需要强调的是，与 CaMN^[5]不同，本文在对抗训练中同样不将片段前 N 帧前置动作输入判别器。前置动作帧属于可观测上下文条件，其内容在训练阶段为真实动作，在推理阶段为历史缓存或上一片段输出；它们并非模型需要生成的目标。因此，本文仅对生成区间 $\hat{\mathbf{g}}_k$ 与 \mathbf{g}_k 进行对抗判别，使对抗目标严格作用于模型实际生成的部分，并与第4.4节的监督损失范围保持一致。

4.5.2 对抗损失定义

判别器的训练目标是区分真实序列与生成序列，其损失定义为：

$$\mathcal{L}_D = -\mathbb{E}_{\mathbf{g}_k} [\log Dis(\mathbf{g}_k)] - \mathbb{E}_{\hat{\mathbf{g}}_k} [\log (1 - Dis(\hat{\mathbf{g}}_k))]. \quad (4.22)$$

生成器则希望其输出被判别器判定为真实，从而对应的对抗损失为：

$$\mathcal{L}_{adv} = -\mathbb{E}_{\hat{\mathbf{g}}_k} [\log Dis(\hat{\mathbf{g}}_k)]. \quad (4.23)$$

其中 $\mathbb{E}_{\mathbf{g}_k}$ 与 $\mathbb{E}_{\hat{\mathbf{g}}_k}$ 分别表示对真实序列与生成序列的采样期望。该对抗目标从分布层面鼓励生成序列在整体运动统计特性上接近真实数据，从而补充监督损失在局部误差上的约束。

4.5.3 交替优化策略

在训练过程中，本文采用交替优化方式更新生成器与判别器。具体而言，对于每个训练批次，我们首先固定生成器参数 θ ，最小化式(4.22)更新判别器参数；随后固定判别器参数，最小化总体损失 \mathcal{L}_{total} （见式(4.15)）更新生成器参数，其中对抗项 \mathcal{L}_{adv} 由式(4.23)给出。通过上述训练方式，判别器不断提升对真实与生成序列的区分能力，而生成器则在监督约束与分布约束的共同作用下，逐步生成更加自然且时间一致的动作序列。

4.6 本章小结

本章围绕严格因果的实时手势生成任务，系统阐述了本文的滑动窗口训练与推理策略，并给出了与之配套的片段级优化目标。首先，我们建立了单步因果预测器作为基本计算单元：采用单向 LSTM 以满足严格因果约束，并在每个预测步中结合跨步循环状态与显式历史动作上下文，采用“短期显式条件 + 长期隐式记忆”的双时间尺度建模，从而提升自回归生成过程的稳定性与长期规律记忆能力。随后，我们沿用 CaMN^[5]的固定长度片段切割方法将长序列组织为训练样本，并在片段内部执行滑动

窗口展开：通过预热阶段缓冲前置动作帧，再在生成阶段逐帧自回归预测并将输出写回历史缓存，最终拼接生成的所有动作帧得到生成序列。

在训练目标方面，本章将监督损失定义在片段生成区间的拼接输出上，采用姿态、速度与加速度的多尺度 Huber 约束以缓解抖动与速度漂移，并进一步引入片段级判别器以提供分布一致性的对抗监督。本文将前置动作帧视为纯条件上下文而非生成目标，在监督损失与对抗训练中均显式移除其影响，使训练目标严格作用于模型实际生成的区间，从而与实时推理阶段的因果生成流程保持一致。

第 5 章 实验结果与分析

为评估 FaceCapGes 的生成质量与实时性能，本文从用户评估、客观指标测量与推理效率三个方面展开实验，并在生成质量上与代表性方法进行对比分析。

5.1 训练配置

本文模型 FaceCapGes 基于 PyTorch 实现，所有实验在单张 NVIDIA RTX 4090 GPU 上进行。

本文基于 BEAT 数据集^[5]进行训练与评估。该数据集包含多模态同步的语音、面部 blendshape 与全身动作信息，以 15 fps 记录多位专业表演者的演讲片段，覆盖多种语义与情绪场景。其标准骨架结构如图 5.1 所示，BEAT 数据集共包含人体中的 47 个关节节点，包括上肢及躯干的三个主要控制点（蓝色区域所示）。下肢关节则保持静态。

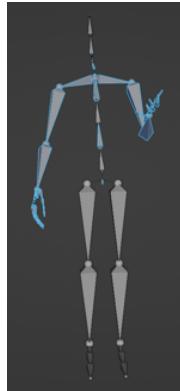


图 5.1 BEAT 数据集的骨架拓扑结构与驱动范围

Figure 5.1 Skeleton Topology and Actuated Joint Range in the BEAT Dataset

本文选取表演者 ID 2、4、6、8 的数据进行训练与测试，其中 2、4 为男性，6、8 为女性，确保在性别与说话风格上的分布均衡。训练集与测试集均包含相同的表演者，但使用不同的演讲片段，在预处理阶段已进行严格划分以避免片段交叉。

训练配置 训练时输入窗口的前序帧数设为 $N = 16$ ，预测长度为 $M = 34$ ，训练片段的切割步长为 10 帧。相邻片段因此存在部分重叠，从而在保证充分上下文信息的同时提升数据覆盖率与时间连续性。批大小设为 256。

优化器采用随机梯度下降 (Stochastic Gradient Descent, SGD)^[37]。基础学习率设置为 $\text{lr}_{base} = 2.5 \times 10^{-4}$, 并根据批大小按线性规则缩放为

$$\text{lr}_g = \text{lr}_{base} \cdot \frac{\text{batch_size}}{128}, \quad (5.1)$$

其中 lr_g 为生成器 (手势生成网络) 的实际学习率。

在对抗训练阶段, 判别器使用相同类型的 SGD 优化器, 其学习率按权重系数 w_d 缩放为

$$\text{lr}_d = w_d \cdot \text{lr}_g, \quad (5.2)$$

本文中设定 $w_d = 0.2$ 。

为防止早期训练阶段的不稳定, 对抗项在第 10 个 epoch 后引入, 即前 10 个 epoch 仅优化重构与时序平滑损失, 从第 11 个 epoch 起加入判别器并交替优化生成器与判别器参数。整体训练共 374 个 epoch。

姿态表示 所有身体动作均转换为连续可微的 Rot6D 表示, 使用 EMAGE^[11] 中的实现方法, 以避免欧拉角奇异性与四元数的符号不确定性。

运行性能 在实时推理阶段, FaceCapGes 以 15 FPS 的速度驱动虚拟角色。

5.2 实验配置

5.2.1 实验对比模型

我们选取 FaceCapGes 的基线模型 CaMN^[5], 以及扩散模型方法中具有代表性的 DiffSHEG^[14] 作为对比模型。

表 5.1 总结了各模型的输入输出模态特征。其中, * 表示模型结构未显式输入说话人 ID, 采用每位说话人独立训练的设置, 说话人 ID 通过模型参数指定。

值得注意的是, FaceCapGes 是唯一满足严格因果约束的在线模型, 因此在评估时采用逐帧推理并将输出拼接为完整序列, 以模拟实时输入流。

5.2.2 跨模型评估设置

本章所有实验均基于 BEAT 数据集进行, 各模型使用一致的骨架拓扑与姿态表示, 从而保证输出格式可直接对齐与比较。其中 CaMN 与 DiffSHEG 使用其官方公开的预训练模型; FaceCapGes 则在相同的数据预处理与骨架设定下由本文训练得到。

表 5.1 对比模型的输入输出模态
Table 5.1 Comparison of Input/Output Modalities of Evaluated Methods

模型	输入模态					未来信息	输出
	音频	面部捕捉	头部姿态	说话人 ID	情绪		
CaMN	✓	✓	✗	✓	✓	✓	身体
DiffSHEG	✓	✗	✗	✓	✗	✓	身体 + 面部
本模型	✓	✓	✓	*	✗	✗	身体

对于离线模型（CaMN 与 DiffSHEG），我们将整段演讲作为输入并一次性生成完整动作序列；对于在线模型（FaceCapGes），我们在关闭批处理（Batch=1）的条件下模拟逐帧输入流，并将逐帧输出按时间顺序拼接得到最终序列，以还原实时交互场景下的运行状态。

5.3 用户评估

为验证模型在交互环境中的表现，本文进行了用户主观评估实验，比较 FaceCapGes、CaMN^[5]、DiffSHEG^[14]三个模型在动作自然性、同步性与多样性方面的主观质量。本节首先介绍用户评估系统与实验配置，随后报告主观评价结果与分析。

5.3.1 用户评估系统与实验配置

实验材料与呈现方式 用户评估所使用的手势动画均基于 BEAT 数据集中的测试集语音片段生成，并以 Biovision Hierarchy（BVH）文件形式保存。BVH 是一种通用的动作捕捉数据格式，通过层级化定义骨骼结构与帧级旋转参数，可直接导入 3D 动画与虚拟人系统。

本文使用的 BVH 文件采用欧拉角旋转表示。由于三种对比模型的输出旋转参数形式不同，因此在导出 BVH 之前，需将输出姿态统一转换为欧拉角表示，以便在后续动画播放中使用统一渲染流程。

本系统当前支持的虚拟人三维模型，需同时具备骨骼绑定，和与 ARKit^[33]兼容的 BlendShape 参数。基于此约束，本实验采用 BEAT 数据集提供的公开的男女两名演讲者三维模型，二者均满足兼容要求，可实现身体与面部的联合驱动。此外，经过 Unity^[38]的 Mecanim^[39]自动骨骼绑定系统匹配，可自动配对模型生成的 BVH 文件中定义的骨骼层级与虚拟人模型的骨骼节点，从而在不依赖手动权重绘制的情况下完成动作重定向。

播放系统实现 我们基于 Unity 自行编写播放脚本，将各模型生成的 BVH 动画用于驱动虚拟人身体骨骼，同时以面部捕捉序列驱动 BlendShape 表情参数，并同步播放原始语音音频。系统支持同时呈现三种模型生成的动画结果：用户可在同一画面中（左、中、右）并行观察三种手势表现，所有语音与面部表情完全一致，唯一变量为身体动作。该设计使参与者能够直接比较不同模型在动作风格、节奏响应与语音同步性方面的差异。

为确保主观评价的公正性与可重复性，系统在每次实验开始前会随机分配三种模型的位置（左、中、右），界面上不会显示模型名称，从而避免潜在偏向。各测试片段的播放顺序在实验前统一设定，以保证不同参与者之间的样本顺序均衡。实验员在播放系统后台记录当前序列与模型对应关系，以便后续结果统计。

实验界面与设备 用户评估系统提供桌面端与虚拟现实（VR）端两种版本，功能完全一致。VR 版本基于 PICO 设备^[40]实现；桌面版支持多窗口并行播放，方便用户同时对比。如图 5.2 所示，播放界面在两种设备上保持统一布局，播放完成后参与者需通过交互界面对三个模型进行排序打分。

VR 用户在沉浸式环境中逐一观看三段动画；桌面端用户则可在单屏上同时观察全部模型。因此前者注重细节感知与临场性，后者更有利于整体风格与节奏的一致性对比。

实验流程与指导 实验正式开始前，研究人员向参与者说明了三项主观评价标准的含义，确保所有被试对评分维度理解一致：

- **真实性：**整体动作是否自然流畅，是否存在明显的违和感，如朝向异常或突然抖动；
- **同步性：**手势动作与语调、语音节奏是否协调一致；
- **多样性：**手势是否丰富多变，避免长时间静止或重复单一动作。

在实验过程中，VR 版本于线下环境进行，桌面版通过线上远程环境执行。两种形式均保持实时交流通道，研究人员可在参与者提问时即时解释操作或澄清评分标准。在正式评估阶段，参与者可多次重播当前片段，但不能返回查看先前内容，以减少记忆偏差。所有播放条件（Unity 场景内的相机角度、光照参数、音量与分辨率设置）在全部被试环境中保持一致，以确保渲染输出的可比性。

需要说明的是，对于 VR 实验，所有测试均在相同的线下实验室环境中进行，使用同一套 PICO 设备与照明条件；而桌面端实验通过远程方式执行，参与者在各自电

脑上运行实验程序。研究人员可通过实时屏幕共享观察其操作流程并保持语音沟通，但无法严格控制其所在房间的光照或环境噪声条件。因此，桌面端实验在观看环境上存在一定差异，但由于任务内容与播放系统完全相同，且实验员在测试中持续指导，可认为该差异对结果的总体影响有限。

实验材料与任务设计 评估样本来自 BEAT 数据集中四位演讲者（ID: 2、4、6、8），其中 2、4 为男性，6、8 为女性。每位演讲者各选取两段平均长度约 1 分钟的语音片段，演讲话题互不重复，共组成 8 段固定视频样本。所有实验均使用相同的 8 段样本，但其呈现顺序在不同被试间经过随机化或平衡化处理，以避免顺序效应。每段视频均包含三种模型生成的动作版本（FaceCapGes、CaMN、DiffSHEG），并在播放时随机分配每个模型的动画在屏幕中的排序。参与者在观看每一片演讲音频后，根据三项主观标准（真实性、同步性、多样性）对三个模型的手势动画表现进行排名评估。

图 5.2 中为用户评估工具的实机界面。画面中共有 3 个虚拟人模型水平分布，在每一片演讲音频播放时，将 3 个生成模型的动画随机分配给 3 个虚拟人的身体骨骼。



图 5.2 用户评估工具实机界面
Figure 5.2 User Study Interface on the Evaluation Device

实验参与者 本实验共邀请 16 名参与者（12 名使用 VR 设备，4 名使用桌面端），涵盖不同性别。所有参与者在实验前均接受了操作说明与校准，并在系统指导下完成评

分练习。

为避免呈现顺序对主观印象造成偏差，另设计了采用平衡拉丁方（Balanced Latin Square）顺序的实验版本，使不同参与者观看样本的顺序均衡分布。该版本实验共招募 8 位 VR 用户（4 男 4 女），排序顺序由 HCI 用户评估工具包^[41]自动生成，确保模型与演讲者组合的呈现顺序在全体被试间均匀分布。所有条件保持一致，唯一变量为视频播放顺序。

实验环境说明 为全面验证模型在不同交互场景下的表现稳定性，本次主观评估设置了桌面端与 VR 端两种实验环境，确保覆盖常规屏幕交互与沉浸式交互两类典型应用场景，具体环境配置如下：

- **桌面端环境：** 参与者通过个人电脑或实验室台式机进行评估，实验界面采用三窗口并行布局，参与者可同时观察左、中、右三个区域的虚拟人动作，聚焦于整体动作风格、节奏同步性的直观对比，注意力分布于整个屏幕的动作全局表现。
- **VR 端环境：** 基于 PICO VR 设备^[40]搭建沉浸式评估场景，参与者佩戴 VR 头显后进入虚拟观测空间，虚拟人以 1:1 比例呈现在眼前，观看距离模拟真实人际交流（约 1.2m）。该环境下参与者注意力更易聚焦于虚拟人上半身细节动作，对空间一致性、动作协调感的感知更敏锐。

两种环境的实验流程、评估指标定义及测试样本完全一致，仅通过设备差异构建不同的观察视角与注意力聚焦模式，以验证模型表现的跨设备适配性。

5.3.2 结果与分析

本节综合分析两轮用户评估的统计结果与参与者反馈。所有结果基于 BEAT 测试集中 4 位演讲者（ID:2、4、6、8；2 男 2 女）各 2 段语音片段，共 8 段固定样本。

总体测评结果 如图 5.3 所示，在 16 名参与者的总体评价中，FaceCapGes 在三个维度（真实性、同步性、多样性）上均优于基线模型 CaMN，并在“真实性”维度上略优于离线扩散模型 DiffSHEG。这一结果表明，FaceCapGes 虽在严格的实时因果约束下运行，但仍能保持与非实时生成模型相近的动作自然度与流畅性。

平衡拉丁方实验结果 图 5.4 展示了平衡拉丁方设置下 8 位 VR 用户的独立结果，该版本严格控制了模型与演讲者组合的呈现顺序。结果与总体趋势一致，且 FaceCapGes

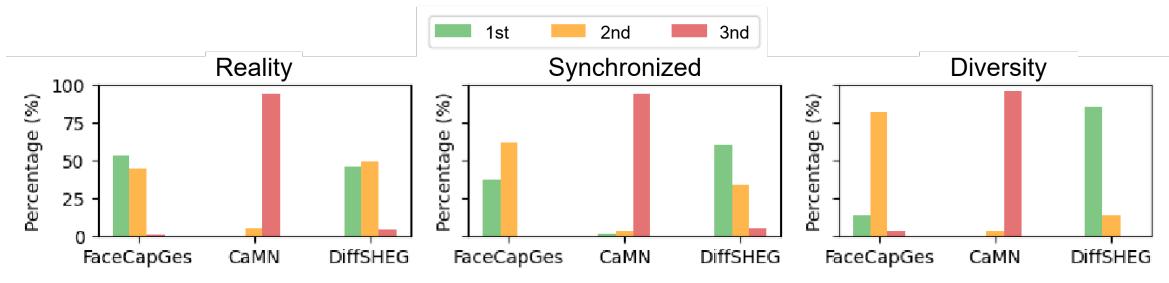


图 5.3 用户评估总体主观排名结果
Figure 5.3 Overall Subjective Ranking Results

在“真实性”与“同步性”得到了更好的评价。这表明实验结论在不同顺序条件下保持稳定，进一步验证了模型主观评价结果的鲁棒性。

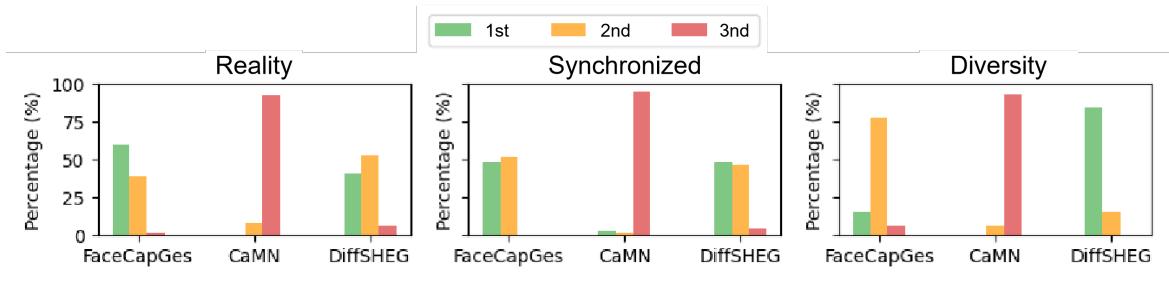


图 5.4 平衡拉丁方设置下的主观排名结果
Figure 5.4 Subjective Ranking Results under Balanced Latin Square Design

用户反馈分析 根据实验后访谈汇总，参与者普遍认为 FaceCapGes 的动作过渡自然、节奏感强，手势响应与语音重音、语调变化更为一致。我们认为，FaceCapGes 融入的头部姿态信息使手势与身体朝向更贴近真实动作，可能是其获得较高真实性评价的重要因素之一。

相比之下，CaMN 在头部与上身动作衔接处常出现僵硬或转向延迟的现象，且头部朝向容易偏离听众方向，从而影响了整体自然度与同步性评分。

对于 DiffSHEG，多数参与者提到其动作丰富度较高，并倾向于将其视为最具多样性的模型。这可能与 DiffSHEG 在生成过程中依赖完整文本输入有关：文本语义边界为动作变化提供了更明确的触发信号，使其更容易生成幅度更大、变化更频繁的动作模式。相比之下，FaceCapGes 主要基于实时语音流与感知模态驱动，缺乏显式语义解析，因此更难实现细粒度的语义级动作响应。

不过，不少参与者也提到 DiffSHEG 在部分片段中存在短暂的手部摆动过快或突然抖动的问题，从而降低了其真实性或同步性的评分。该现象可能与 Axis-Angle 旋

转表示在数值空间中存在非连续点或奇异性有关，从而在训练或推理过程中更容易诱发关节角度的突变与抖动。相比之下，本文模型采用 Rot6D^[10] 旋转表示，该表示具有更好的连续性与数值稳定性，有助于缓解由旋转表示引入的突变问题。在本次用户研究中，FaceCapGes 未出现明显的抖动相关反馈，生成动作在关节转动与姿态过渡上保持较为连贯自然。

结果讨论 FaceCapGes 的因果式时间建模与头部姿态融合策略有效提升了局部动作的平滑性与节奏协调，在不依赖未来输入信息的条件下完成逐帧推理，更契合在线实时交互场景的因果性约束。此外，平衡拉丁方版本进一步证明主观结论在不同呈现顺序下的一致性，排除了顺序偏差对结果的显著影响。

综上，用户研究表明 FaceCapGes 在在线实时生成条件下仍能维持与离线模型相近的主观表现，验证了本文提出的多模态融合与时间建模策略的有效性。

5.4 定性分析

我们强烈建议观看附录中的演示视频（附录 A），内容包含 Ground Truth (GT)、本模型 (FaceCapGes)、CaMN 与 DiffSHEG 在不同演讲数据上生成的手势的并排展示，能够直观体现时间对齐性、手势响应性以及头部-身体协调性方面的差异。

为进一步理解模型在实时因果约束下的动态响应行为，我们对每个模型的生成手势中进行了逐帧观察，并选取具有代表性的片段进行案例分析。

5.4.1 生成动作平滑性

如图 5.5 所示，FaceCapGes 在多个片段中均能平滑地响应说话人的语调变化。当语调出现明显上升或下降趋势时，本模型生成的双手高度能够随之连续变化，且动作幅度保持自然，整体运动轨迹连续、关节过渡平稳。该现象与主观用户评估中参与者对 FaceCapGes “动作过渡自然、节奏感强”的反馈一致。

相比之下，CaMN 的生成动作在语调变化较快的片段中表现出一定的迟滞性：双手高度调整通常较为缓慢，且动作幅度变化更趋于保守，导致整体动作轨迹的动态范围较小，视觉上更容易产生“僵硬感”。这一观察与用户反馈中提到的 CaMN “转向延迟与衔接僵硬”现象相一致。

DiffSHEG 的生成动作整体更活跃，在部分片段中能够在更细粒度的时间尺度上产生频繁的手势变化。然而，我们也观察到其在少数时间步出现局部关节的突然加速

或短暂抖动，表现为手部轨迹在相邻帧间产生明显跳变或方向快速反转。该现象与主观用户反馈中提到的“偶发抖动”一致，可能与其旋转表示在数值空间中存在非连续点有关，从而使推理阶段更容易产生局部突变。

总体而言，FaceCapGes 通过采用 Rot6D 表示并结合帧级因果时间建模，在保证实时性的同时维持了更高的动作连续性与稳定性，生成结果在视觉平滑性上更接近真实动作序列。

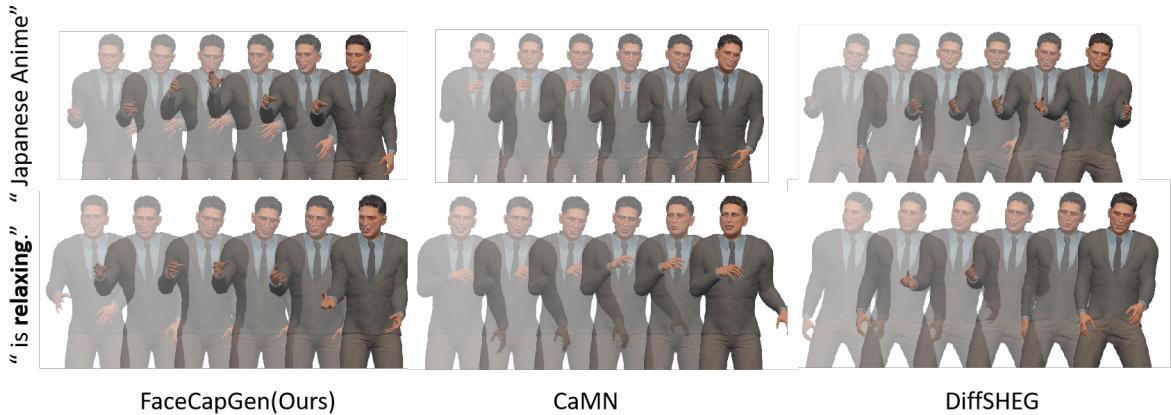


图 5.5 生成动作效果对比
Figure 5.5 Qualitative Comparison of Generated Gestures

5.4.2 头部朝向与手势空间指向一致性

为分析头部姿态输入对生成手势空间指向的影响，我们选取测试集中 GT 头部朝向发生明显偏转的片段，并对比不同模型在该时刻生成动作的朝向一致性表现。这里的“空间指向一致性”指生成手势的主要运动方向是否与角色的头部/身体朝向保持匹配，从而反映模型是否能够利用非语言姿态线索生成更符合交互场景的空间表达。

值得注意的是，CaMN 与 DiffSHEG 在输入端均不显式包含头部姿态信息，因此其生成的手势方向主要依赖语音或文本内容，难以体现真实头部转向带来的空间指向变化。

图 5.6 展示了一个包含明显头部转向变化的片段截图，从左到右依次为 GT、本模型（FaceCapGes）、CaMN 与 DiffSHEG 的生成结果。在该片段中，GT 的身体与头部朝向呈现出明确的空间指向，并且手势多沿着当前头部/躯干朝向展开，体现出面向不同听众时的自然交流习惯。

我们观察到，本模型在多数帧中能够保持与 GT 相近的身体朝向，并使双手动作的空间指向与头部朝向保持一致，例如在头部向右侧转动的阶段，生成的手势也倾向



图 5.6 头部朝向与手势指向一致性的动作对比

Figure 5.6 Comparison of Head Orientation and Gesture Direction Consistency

于朝向相同方向展开。相比之下，CaMN 的动作表现更为僵硬，身体朝向变化较弱且缺乏稳定的空间锚定，手势方向在多帧间呈现随机波动。DiffSHEG 虽能在其可获取

的语音信息下生成较连贯的手势节奏与整体动作幅度，但由于其输入不包含头部姿态信号，生成结果难以反映 GT 中由头部转向引起的空间指向变化。

该现象表明，在严格因果与实时输入条件下，头部姿态作为额外的非语言模态能够为手势生成提供空间指向上的中层约束，使动作更自然地与说话者的注意方向和交互对象匹配。这种头-手空间一致性在元宇宙等虚拟交互场景中尤为重要：当用户面向不同方向的听众或对象进行交流时，头部转向与与之匹配的手势能够增强空间合理性与沉浸感，从而提升多方位交流中的可理解性与聆听体验。

5.5 客观评估指标与实现细节

为客观层面评价模型在动作自然性、节奏同步性与多样性等方面的表现，本文在客观评估中采用四项度量：Fréchet 手势距离（Fréchet Gesture Distance, FGD）^[42]、语义相关动作召回率（Semantic Relevance Gesture Recall, SRGR）^[5]、节奏对齐度（Beat Alignment, BA）^[5,43]以及 L1 范数（L1DIV）。这些指标分别对应生成动作在分布一致性、语音同步性与变化丰富性等不同维度，共同构成对模型质量的综合评估体系。

5.5.1 Fréchet 手势距离（Fréchet Gesture Distance）

FGD^[42]用于衡量生成手势分布与真实手势分布之间的统计距离，灵感源自图像生成领域的 Fréchet Inception Distance (FID)。不同于图像任务直接利用 Inception 网络特征，在动作生成领域，特征空间需由单独训练的动作自编码器定义。该自编码器通过重构任务学习手势的潜在表示，使潜在空间具备对运动模式的压缩与区分能力。在该潜在空间中，假设真实分布与生成分布的高维嵌入向量分别为 $\mathcal{N}(\mu_r, \Sigma_r)$ 与 $\mathcal{N}(\mu_g, \Sigma_g)$ ，其中， $\mathcal{N}(\cdot)$ 为高斯分布， μ 与 Σ 分别为嵌入特征的均值向量与协方差矩阵。

此时，FGD 定义为：

$$\text{FGD} = \|\mu_r - \mu_g\|_2^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}). \quad (5.3)$$

其中， $\text{Tr}(\cdot)$ 为矩阵迹运算。

较小的 FGD 值表示生成动作的统计分布更接近真实数据，可反映动作的整体自然度与风格一致性。

FGD 的模型结构与训练配置 本文在每位说话人的训练集上分别训练一组评估用自编码器，以避免跨说话人分布差异对指标的干扰。自编码器输入为以 Rot6D 表示的

上半身骨架序列，训练目标为最小化位置与速度的重构误差：

$$\mathcal{L}_{AE} = \lambda_r \|\hat{\mathbf{g}} - \mathbf{g}\|_2^2 + \lambda_v \|\hat{\mathbf{g}}' - \mathbf{g}'\|_2^2, \quad (5.4)$$

其中 $\lambda_r = 1, \lambda_v = 0.1$, \mathbf{g}' 为速度序列。

训练配置如下：输入片段长度设为 32 帧，训练片段切割步长设为 10 帧，批大小设为 256，编码器与解码器的隐藏层维度均为 128。优化器采用 SGD^[37]，基础学习率设为 $lr_{base} = 1.2 \times 10^{-4}$ ，并按批大小线性缩放为

$$lr = lr_{base} \cdot \frac{\text{batch_size}}{128}. \quad (5.5)$$

训练过程中仅使用位置与速度重构损失（式（5.4）），不包含加速度或对抗项。自编码器共训练 400 个 epoch。

用于 FGD 评估的骨架感知自编码器 在基线模型 CaMN 的 FGD 评估中，采用了一种基于时间卷积的嵌入式自编码器（embedding-based autoencoder）进行手势特征提取。该结构将每一帧姿态平铺为高维向量输入，并对各关节分量进行独立建模。在实际使用中我们观察到，这类嵌入空间对旋转表示的数值尺度较为敏感。当采用 Rot6D 表示时，不同关节与分量之间的不均衡方差可能在潜在空间中被进一步放大，从而导致协方差估计条件较差，并引发 FGD 数值不稳定的问题。

为提高 FGD 评估的鲁棒性，本文采用了一种骨架拓扑感知的自编码器作为特征提取器。该模型在编码阶段显式引入骨架邻接矩阵 A ，并通过对相邻关节进行局部卷积实现结构约束。具体而言，第 l 层中关节 i 的特征向量 $\mathbf{h}_i^{(l)} \in \mathbb{R}^{d_l}$ ，其中 d_l 表示第 l 层中每个关节节点的特征维度。通过聚合其邻域 $N(i)$ 内相邻关节 j 的特征 $\mathbf{h}_j^{(l)}$ 进行更新，其中 A_{ij} 表示骨架邻接矩阵中节点 j 到节点 i 的连接权重，即

$$A_{ij} = \begin{cases} 1, & \text{若关节 } i \text{ 与 } j \text{ 在骨架拓扑中直接相连, 或 } i = j; \\ 0, & \text{否则.} \end{cases} \quad (5.6)$$

此外，设 $W^{(l)} \in \mathbb{R}^{d_{l+1} \times d_l}$ 与 $b^{(l)} \in \mathbb{R}^{d_{l+1}}$ 分别为第 l 层的可学习权重矩阵与偏置项。设 $\sigma(\cdot)$ 为非线性激活函数，在本文实现中取为双曲正切函数 $\tanh(\cdot)$ 。上述局部邻域聚合过程可形式化表示为：

$$\mathbf{h}_i^{(l+1)} = \sigma \left(\sum_{j=1}^J A_{ij} W^{(l)} \mathbf{h}_j^{(l)} + b^{(l)} \right), \quad (5.7)$$

这种基于局部邻域的权重共享机制有助于保持人体运动的空间结构一致性，并在特征提取过程中缓解由旋转表示差异带来的数值尺度放大问题。

实验结果表明，该骨架感知自编码器在 Rot6D 下产生更加稳定的潜在分布统计量。在本文的实验设置中，该设计有效提升了 FGD 评估的数值稳定性，并有效避免了在 Rot6D 表示条件下使用嵌入式自编码器时出现的极端 FGD 数值现象。

5.5.2 语义相关动作召回率 (Semantic Relevance Gesture Recall)

SRGR 指标 (Semantic Relevance Gesture Recall)^[5] 用于衡量生成手势在语义相关时间段内与真实手势在数值层面的匹配程度。该指标关注的是语音语义触发的关键手势是否被正确生成，而非整体分布一致性或语音-手势节奏相关性。

与基于分布的评估指标（如 FGD）不同，SRGR 属于基于阈值的逐帧召回率度量，通过统计生成手势在允许误差范围内命中的比例，反映模型在语义相关动作重现方面的准确性。

定义与原理 在本文实现中，SRGR 作用于关节的旋转表示。设真实手势序列与生成手势序列在第 t 帧第 j 个关节的旋转表示分别为 $\mathbf{r}_{t,j}$ 与 $\hat{\mathbf{r}}_{t,j}$ ，其中 $\mathbf{r}_{t,j} \in \mathbb{R}^6$ 为关节的 Rot6D 表示， T 为序列总帧数， J 为关节数量。

在给定旋转表示误差阈值 δ 的条件下，若生成关节旋转与真实关节旋转之间的表示差异满足 $\|\mathbf{r}_{t,j} - \hat{\mathbf{r}}_{t,j}\|_1 < \delta$ ，则认为该关节在该时刻被成功召回。在本文实验中，阈值固定设为 $\delta = 0.5$ 。

SRGR 通过对所有时间帧与关节进行统计，并引入语义相关性权重 λ_t ，定义为：

$$\text{SRGR} = \frac{1}{TJ} \sum_{t=1}^T \sum_{j=1}^J \lambda_t \mathbb{I}(\|\mathbf{r}_{t,j} - \hat{\mathbf{r}}_{t,j}\| < \delta), \quad (5.8)$$

其中 $\mathbb{I}(\cdot)$ 为指示函数。

语义相关性权重 λ_t 用于强调语音中语义显著时间片段（如强调词或语义触发点）对应的手势匹配程度，由 BEAT 数据集提供的语义标注确定，从而使 SRGR 更加关注语义相关手势的召回情况，而非对所有时间片段进行均匀统计。

5.5.3 节奏对齐度 (Beat Alignment)

Beat Alignment (BA) 指标用于衡量语音节拍事件与手势关键动作事件在时间轴上的对齐程度，反映模型在语音-手势时序同步性 (temporal synchronization) 方面的性能。本文采用 BEAT/CaMN 中使用的 BeatAlign 指标^[5]，其原始形式由 Li 等人

提出于 AI Choreographer 工作中^[43]，并可视为音频节拍与动作节拍集合之间的单向 Chamfer 相似度度量。

与 SRGR 基于阈值的逐帧数值匹配不同，BA 关注的是离散事件层面的时间对齐关系，即手势关键动作是否在时间上合理地响应了语音中的重读节拍。

定义与原理 设语音节拍事件集合为 $\mathcal{B}_s = \{t_k^s\}$ ，通过对语音信号的能量相关特征进行起始点检测（onset detection）得到。具体而言，首先基于谱能量变化计算 onset strength 曲线，并在该曲线上进行峰值检测以获得候选语音事件位置。随后，引入 Root Mean Square (RMS) 特征作为短时能量幅度的描述，并在 RMS 曲线上对检测到的起始点进行 backtracking 校正，从而将语音节拍事件定位至能量实际开始上升的位置，以提高时间定位的稳定性与准确性。

手势关键动作事件集合为 $\mathcal{B}_g = \{t_m^g\}$ ，定义为关节运动速度的局部极小值点，对应动作中的停顿或方向变化等显著运动事件。

对于每一个语音节拍事件 t_k^s ，计算其与所有手势关键动作事件之间的最小时间偏差：

$$\Delta t_k = \min_m |t_k^s - t_m^g|. \quad (5.9)$$

随后采用高斯核函数将时间偏差映射为相似度分数，从而得到单个语音节拍的对齐得分。最终 BA 指标定义为：

$$BA = \frac{1}{|\mathcal{B}_s|} \sum_k \exp \left(-\frac{(\Delta t_k)^2}{2\sigma^2} \right), \quad (5.10)$$

其中 σ 为时间尺度参数，本文中取 $\sigma = 0.3$ ，用于控制对齐容忍范围。

该定义可视为从语音节拍集合到手势节拍集合的单向 Chamfer 相似度，当语音节拍与手势关键动作在时间上高度对齐时，BA 值接近 1；反之，当二者时间偏差较大时，BA 值趋近于 0。

本文使用 BA 指标评估模型在语音重读节拍与手势关键动作之间的时间同步能力，该指标与主观观察到的语音-手势同步自然度通常具有较好一致性。

5.5.4 L1 范数

L1 范数 (L1DIV) 指标^[5]用于衡量模型生成手势序列的多样性，即不同生成样本之间在动作空间中的差异程度。该指标反映模型在避免生成结果收敛到平均动作模式 (mode collapse) 的同时，是否能够保持足够丰富的动作变化。

定义与原理 设模型在评估过程中生成 N 个手势序列样本，第 i 个生成样本在第 t 帧第 j 个关节的旋转表示为 $\hat{\mathbf{r}}_{t,j}^{(i)}$ ，其中 $\hat{\mathbf{r}}_{t,j}^{(i)} \in \mathbb{R}^6$ 为关节的 Rot6D 表示， T 为序列总帧数， J 为关节数量。

L1DIV 通过计算不同生成样本之间在所有时间帧与关节上的平均 L1 距离，来刻画生成动作分布的离散程度。其数学形式定义为：

$$\text{L1DIV} = \frac{1}{N(N-1)} \sum_{i < k} \frac{1}{TJ} \sum_{t=1}^T \sum_{j=1}^J \left\| \hat{\mathbf{r}}_{t,j}^{(i)} - \hat{\mathbf{r}}_{t,j}^{(k)} \right\|_1. \quad (5.11)$$

较高的 L1 范数表明模型生成具有较强的多样性，但过高可能意味着动作不稳定或噪声放大。因此，L1 范数通常与 FGD 联合分析：FGD 反映真实度，L1 范数反映丰富度，两者共同平衡模型在自然性—多样性维度上的表现。

5.5.5 评估区域设定与公平性说明

由于本文方法 FaceCapGes 在输入端显式引入了真实的头部姿态作为额外模态，并在输出中生成包含头部旋转的上半身骨骼序列，若直接将头部旋转也计入整体误差，可能会对不使用头部姿态输入的对比方法（如 CaMN 与 DiffSHEG）造成不公平的优势。为保证跨方法的可比性，本文在所有定量评估指标中分别报告两种评估区域：

- (1) 上半身（含头部）：包含头部与上半身全部关节旋转；
- (2) 上半身（不含头部）：在计算指标时移除头部关节的旋转分量，仅统计身体与手臂部分。

具体而言，在计算 FGD、SRGR、BA 与 L1DIV 时，我们将头部关节的旋转维度从指标计算中排除，从而消除头部旋转差异对指标的直接影响。

5.6 定量评估结果

表 5.2 显示，FaceCapGes 在所有指标上均优于 CaMN。值得注意的是，“上半身（含头部）”与“上半身（不含头部）”两种评估区域下的结果趋势基本一致，表明本文方法的性能优势并非仅来源于头部姿态输入或头部输出的额外信息。

此外，表 5.2 还显示，与 DiffSHEG 相比，本模型 FGD 更低，SRGR 相近，且在 L1 范数上表现最优，表明其生成动作具备良好多样性。但这一结论与用户主观评分存在一定出入：DiffSHEG 在多样性上主观排名更高。

这个现象可能来自于 L1 范数的局限性：它主要衡量空间偏离程度，不能直接体

表 5.2 定量评估结果
Table 5.2 Quantitative Evaluation Results

区域	方法	FGD↓	SRGR↑	BA↑	L1DIV↑
上半身（含头部）	CaMN	47.732	0.098	0.845	7.591
	DiffSHEG	<u>26.846</u>	0.109	<u>0.883</u>	<u>11.028</u>
	本模型	23.385	<u>0.107</u>	0.913	13.284
上半身（不含头部）	CaMN	49.437	0.103	0.845	7.327
	DiffSHEG	<u>26.847</u>	0.110	<u>0.883</u>	<u>10.990</u>
	本模型	23.384	<u>0.107</u>	0.913	13.330

现动作的颗粒度，或用户感知上的丰富性。虽然我们的模型在动作结构上更丰富，但用户普遍认为 DiffSHEG 的动作更活跃，在合理的时机做出了更多吸引注意的动作。

5.7 消融实验分析

我们围绕两个核心设计展开消融：

- (1) 头部姿态输入是否能提供有效空间与时序线索；
- (2) 帧级自回归生成与滑动窗口训练是否优于片段级一次性解码。

为此，我们构造了三种变体：移除头部姿态输入、移除帧级生成策略，以及基线 CaMN。

所有消融实验均基于 BEAT 数据集的第 2 位说话人进行，训练与测试划分与主实验保持一致。

表 5.3 展示了各模块在四个指标上的结果。

表 5.3 消融实验结果
Table 5.3 Ablation Study Results

区域	变体	FGD↓	SRGR↑	BA↑	L1DIV↑
上半身（含头部）	CaMN	32.870	0.111	0.858	7.214
	移除头部姿态	19.591	0.123	0.916	10.642
	移除帧级生成	21.592	0.125	0.892	10.456
	FaceCapGes（本模型）	19.290	0.123	0.918	10.871

头部姿态输入的贡献 对比“移除头部姿态”与完整模型可以发现，尽管两者在 SRGR 上差异不大（均为 0.123），但完整模型在 FGD、BA 与 L1DIV 上均取得更优表现，尤其在动作分布一致性（FGD）与多样性（L1DIV）上呈现出稳定收益。这表明头部姿

态作为非语言模态提供了额外的空间方向与交互焦点线索，能够帮助模型在生成手势时保持更一致的身体朝向与动作指向，从而提升动作自然度与表达丰富性。

需要注意的是，该提升幅度相对温和，原因可能在于语音与面部表情已包含较强的节奏与情绪提示，头部姿态主要补充空间层面的约束，因此其收益更集中地反映在分布与多样性相关指标上。

帧级自回归生成的优势 值得注意的是，“移除帧级生成”版本允许利用未来上下文并采用双向 LSTM，但其整体表现仍不及帧级版本，尤其在 FGD 上出现较为明显的退化（FGD 从 19.290 上升至 21.592）。

一种合理解释是，该变体在训练阶段以独立窗口为单位进行优化，而在测试阶段采用整段演讲作为长序列输入并一次性生成完整动作序列。这种训练与推理的序列长度范围不一致，可能导致模型在长序列推理时的隐状态传播方式偏离训练分布，从而使生成动作在嵌入空间中的统计特性偏离真实数据分布，表现为生成分布与真实动作分布的距离增大（FGD 上升）。

尽管本模型同样采用窗口化训练，但在每个窗口内通过滑动窗口展开与纯自回归预测进行多步生成，并对每步输出的误差累计取平均作为优化目标，使模型在训练阶段即暴露于自身预测分布并学习局部动力学的稳定性。因此训练目标与在线推理过程保持一致，对推理阶段序列长度变化所引入的分布偏移具有更强鲁棒性。

与基线模型 CaMN 的差异 相比基线 CaMN，本模型在所有指标上均显著提升，其中 FGD 从 32.870 降至 19.290，表明生成分布更接近真实动作；同时 BA 与 L1DIV 的提高说明动作更平衡且更具表达多样性。这进一步验证了本文引入的因果时序建模、头部模态补充与滑动窗口训练策略对于实时交互场景下的手势生成具有有效作用。

此外，基线 CaMN 采用欧拉角而本模型采用 Rot6D 表示，该表示差异亦可能部分解释性能提升幅度较大的原因。

5.8 性能评估

5.8.1 单帧推理性能

FaceCapGes 作为端到端实时手势生成框架的核心计算模块，其性能评估聚焦于单帧输入，单帧输出的核心推理流程，即模型接收当前帧的语音、面部表情与头部姿态多模态输入，实时输出对应帧的上半身手势骨骼动画。

为模拟真实应用中的实时输入流场景，性能测试基于 BEAT 数据集的测试集展开，关闭批处理机制，确保每帧数据均独立输入模型进行推理，还原逐帧处理的实际运行状态。测试过程中，我们将测试集中总计 93015 帧的多模态输入数据传入模型，记录从首帧输入到末帧输出的总推理耗时，通过总时长与测试帧数的比值计算平均单帧处理时间。测试时使用的硬件配置为单张 RTX4090 硬件。

测试结果如表 5.4 所示，模型平均单帧处理时间为 6.07 毫秒，具备良好的实时响应能力。

表 5.4 推理速度评估结果

Table 5.4 Inference Speed Evaluation Results

指标	时间
测试帧数	93015 (f)
推理总时长	504 (s)
平均单帧时间	6.07E-03 (s/f)

5.8.2 端到端计算链路延迟

结合第 2.5 节所述的端到端框架流程，系统端到端延迟可按时间顺序拆解为数据采集、特征提取、模型推理与结果返回四个核心阶段。各阶段性能消耗及瓶颈分析如下：

- **数据采集阶段：**依赖设备端传感器实时捕获多模态信号，主要耗时来源于 ARKit 面部与头部姿态追踪。根据官方文档标注，在 iOS 设备上 ARKit 的目标追踪帧率为 60 FPS^[2]，对应输入更新周期约为 16.7 ms。该阶段耗时由设备端硬件算力与系统负载决定。
- **特征提取阶段：**将原始传感器数据转换为模型可识别的结构化特征（语音 Mel 频谱、面部 BlendShape 系数、头部 Rot6D 旋转参数）。测试原型中采用 Librosa 对测试集音频进行离线特征提取，93015 帧数据的总计算耗时约 61 秒，对应的平均单帧计算成本约为 0.66 ms。需要注意的是，该统计不包含实时 I/O 与缓冲管理等系统开销，但可用于估计音频特征提取的计算量级。实际部署时可替换为 PyAudio 等实时流式提取工具，因此该阶段预计不构成主要性能瓶颈。
- **模型推理阶段：**为端到端流程的核心计算环节，结合第 5.8.1 节的模型推理性能测试，在单张 RTX4090 硬件、无批处理（Batch=1）配置下，单帧推理耗时为 6.07 ms。该阶段耗时与硬件算力强相关，是由神经网络结构设计决定的主要性能变量。

- **结果返回阶段：**将模型输出的骨骼姿态参数传输至渲染引擎，耗时可忽略（通常低于 0.1 ms），不构成性能瓶颈。

在本地推理设置下，若将一次响应链路定义为单帧采集完成后立即进入后续计算，则计算链路延迟可近似表示为：

$$t_{\text{e2e}} \approx t_{\text{ARKit}} + t_{\text{feat}} + t_{\text{infer}} + t_{\text{return}}, \quad (5.12)$$

其中 $t_{\text{ARKit}} \approx 16.7 \text{ ms}$, $t_{\text{feat}} < 1 \text{ ms}$, $t_{\text{infer}} = 6.07 \text{ ms}$, $t_{\text{return}} < 0.1 \text{ ms}$ 。因此在该设置下，计算链路理论延迟可视为 25 ms 以下。

需要注意的是，该估计未计入音频特征提取的窗口缓存与渲染端同步可能引入的额外等待，实际交互延迟还将受渲染刷新周期影响。

远程推理部署的额外开销 在移动端等算力受限的平台上，模型推理可部署于远程服务器并通过网络进行输入输出传输。此时端到端延迟需进一步加上网络往返与序列化开销：

$$t_{\text{e2e}}^{\text{remote}} \approx t_{\text{ARKit}} + t_{\text{feat}} + t_{\text{net}} + t_{\text{infer}} + t_{\text{return}}, \quad (5.13)$$

其中 t_{net} 表示网络传输与通信开销，其大小由网络条件与系统实现决定，本文不在此展开测量。

5.8.3 系统更新率

在本实验设置下，模型在单张 RTX4090 上的单帧推理时间为 6.07 ms，对应的理论推理吞吐约为 165 FPS，说明模型推理阶段在当前硬件条件下不会成为实时交互中的性能瓶颈。需要强调的是，该数值仅反映神经网络推理吞吐能力，并不等同于系统端到端更新率。

然而，端到端系统的实际更新率仍将受到输入采集频率与渲染刷新率的共同约束。根据官方文档，ARKit 面部与头部追踪通常以 60 FPS 为目标更新率^[2]，因此在采用 ARKit 作为面部与头部输入来源的应用场景中，系统可获得的相关输入更新频率理论上可视为 60 Hz，相应地，系统的有效输出更新率亦不会超过 60 FPS。在远程推理部署下，网络传输与同步开销可能进一步降低实际更新率。

5.8.4 帧率设定的可扩展性

本文实现中采用 15 FPS 的动作时间采样率进行训练与输出，主要原因在于与 BEAT 数据集预处理及基线模型的设置保持一致，以便进行公平对比。该采样率选择

并不构成本模型的固有限制。

但需要注意的是，本文模型的输入模态不包含帧间时间间隔 Δt 的显式信息，因此训练过程中隐式假设固定的离散时间步长（即 $\Delta t = 1/15\text{ s}$ ）。当部署环境的输入采样率或输出刷新率发生变化（例如 ARKit 为 60 FPS）时，若直接使用 15 FPS 训练的模型进行逐帧推理，模型可能会对运动速度与节奏尺度产生偏差，从而影响手势动态表现。

因此，在目标运行帧率与训练帧率不一致的情况下，更稳妥的做法是对数据集进行对应帧率的重采样并重新训练模型，以确保训练时域与实际系统时域一致，从而使模型学习到正确的动力学时间尺度。未来也可进一步探索将 Δt 作为显式条件输入，提高模型对不同采样率输入的鲁棒性。

在实际部署中，系统可依据端到端管线中的主要瓶颈选择合适的目标动作帧率：当推理计算为瓶颈（如移动端或低算力设备）时，可维持较低帧率以确保每帧按时生成；当输入采样率与计算资源允许（例如更高刷新率的追踪）时，可采用更高帧率训练版本以获得更细粒度的时间响应与动作细节表达。

5.9 本章小结

本章围绕 FaceCapGes 的生成质量与实时性能开展了系统评估。

在主观评估方面，本章结合两轮用户评估对模型在真实感、同步性与多样性三个维度进行了分析，结果显示本方法整体优于基线模型 CaMN，并在真实感维度上与离线模型表现相近。结合用户反馈与访谈，本章进一步讨论了头部姿态信息在提升动作朝向一致性与真实感方面的潜在贡献。此外，通过对生成结果的定性分析，我们观察到本方法能够随头部朝向变化生成方向一致的手势动作，从而增强交互场景下的空间合理性。

在客观评估方面，本文采用 FGD、SRGR、BA 与 L1DIV 等指标，分别从动作分布一致性、语音节奏对齐、动作平衡性与多样性等角度对比了 CaMN 与 DiffSHEG 等方法，并通过“含头部/不含头部”两种评估区域设定保证指标对比的公平性。在此基础上，本章进一步通过消融实验分析了头部姿态输入与帧级自回归生成策略对性能提升的贡献，验证了各模块设计的有效性。

最后，本章对模型推理速度、端到端计算链路延迟与系统更新率进行了评估，表明 FaceCapGes 在保持较高生成质量的同时具备实时交互所需的低延迟响应能力。

第6章 结论

6.1 本文工作总结

本论文围绕“面向用户交互的在线实时数字人驱动”这一核心目标，针对现有语音驱动手势生成方法普遍依赖整句输入、难以满足严格实时约束的问题，提出了一种仅使用在线可获取信号（语音、面部表情与头部姿态）即可逐帧生成上半身骨骼动作的帧级生成方法 FaceCapGes，实现了无需用户实际做出手势的低门槛数字人自然表达能力。论文的主要工作与贡献可总结如下：

(1) 构建了面向在线实时数字人驱动的多模态帧级手势生成系统框架。针对在线交互场景下语音逐字输入带来的未来信息不可用问题，本文从系统层面提出了严格因果的实时生成机制，明确了从信号采集、特征同步、帧级推理到虚拟人驱动渲染的完整流水线结构，并在实时约束下制定了输入模态选择、姿态表示方式以及端到端数据流组织策略。该框架保证了模型能够在仅依赖当前与历史信息的情况下持续输出动作流，为后续的可部署实现提供了系统基础与设计规范。

(2) 在级联多模态架构中引入头部姿态作为新的实时输入模态，并提出弱耦合融合策略以增强节奏与指向一致性。本文在继承 CaMN 级联设计思想的基础上，将头部姿态视为实时可获取的辅助输入信号，用于提供节奏前瞻与空间锚定线索，从而弥补语音模态在严格因果条件下对方向一致性与互动焦点建模不足的问题。为此，论文设计了头部姿态编码器，并通过弱耦合方式将其作为独立通道拼接进多模态隐向量，避免其与语音/表情编码产生过强耦合导致收敛困难。该设计在保持在线实时推理能力的同时，使模型能够显式利用用户头部朝向变化，从而增强生成动作的空间协调性与指向一致性。

(3) 提出并实现了基于滑动窗口展开的自回归训练策略，使训练过程与在线推理严格一致，从而提升实时生成的稳定性与连续性。为解决帧级自回归生成中常见的漂移与抖动问题，本文提出了结合片段切割与滑动窗口展开的训练流程，将单步因果预测器在片段内部逐帧展开，并采用历史动作缓存写回机制以模拟真实在线推理过程，使训练与部署阶段保持一致。此外，论文在片段生成区间内引入姿态、速度与加速度的多尺度监督约束，并结合片段级判别器进行对抗训练，从局部运动统计分布层面进一步提升生成动作的自然性与时间连续性。该训练策略在严格因果约束下有效缓解了自回归误差累积带来的稳定性下降问题。

(4) 构建了完整的评估平台并进行了主客观与性能实验验证, 证明 FaceCapGes 在严格因果条件下仍具备良好的生成质量与实时推理能力。论文搭建了统一的跨模型评估与渲染平台, 使不同方法可在相同输入条件与渲染设置下公平对比, 并在 BEAT 数据集上对 FaceCapGes、CaMN 与 DiffSHEG 等代表性方法开展了用户主观评估、客观指标测量, 并进一步对 FaceCapGes 进行了消融分析与实时性能测试。实验结果表明, 在严格因果约束下, FaceCapGes 的生成真实性可达到与扩散模型方法相当的水平; 同时在韵律变化较强的语音片段及头部方向变化明显的场景中, 其动作表现更平滑且空间指向与真实数据更一致, 体现出该方法在虚拟交互场景中保持空间表达一致性的优势。此外, 模型帧级推理效率与端到端链路延迟满足实时交互应用的运行需求, 验证了其作为可部署在线数字人驱动方案的可行性。

综上所述, 本论文在严格因果约束下, 提出了基于语音、面部捕捉与头部姿态的在线帧级手势生成方法 FaceCapGes, 并从系统设计、模型结构、训练策略与实验验证四个方面证明了在无需未来信息的条件下实现自然随语手势生成的可行性与应用价值。FaceCapGes 在严格实时约束下实现了无需手势采集的自然随语动作生成, 可支持实时虚拟人直播、沉浸式虚拟社交互动等交互式数字人应用场景, 提升了数字人表达的自然性与互动一致性。与此同时, 本研究也为未来进一步融合高层语义信息与预测性控制机制的实时数字人驱动方法提供了技术基础与参考。

6.2 未来工作展望

6.2.1 高层语义信息

当前模型主要关注语音声学特征与运动感知模态对手势生成的影响, 尚未显式引入语言层面的语义理解与表达意图建模。未来可结合实时语音识别与增量式语义解析技术, 引入语篇结构、强调意图或对话功能等高层语义信息, 以丰富手势在交互场景中的表达能力。在不破坏实时性的前提下, 探索对有延迟但可修正的语义假设的鲁棒利用方式, 将有助于提升生成手势在语义层面的准确性与一致性。

6.2.2 面向未来趋势的预测性训练目标

从建模目标的角度来看, 当前 FaceCapGes 的训练过程主要以当前时间步手势姿态的重建误差为优化目标, 即在给定历史与当前多模态输入的条件下, 监督模型对当前手势的预测精度。然而, 该学习目标并未对未来时间段内手势节奏与结构变化施加显式约束, 使模型对历史信息的利用更多服务于当前帧生成, 而非对即将发生动作

变化进行前瞻性建模。

未来的研究可在现有框架基础上，引入针对手势未来趋势的预测性监督信号，尤其是充分挖掘头部与面部动态中所蕴含的准备性线索。与直接预测未来具体手势姿态不同，该方向更侧重于对抽象化时序属性的建模，例如未来短时间窗口内的手势起始概率、运动能量变化或强调强度等。这类趋势性变量具有时间平滑、语义明确且可提前出现的特点，适合作为实时系统中的前瞻性约束。

通过在训练阶段同时优化当前手势生成与未来趋势预测两个目标，模型有望学习到更具时间结构性的中间表示，从而在不引入额外模态或显著增加系统延迟的前提下，实现对手势节奏的提前准备与更稳定的时序对齐。

参考文献

- [1] KARTYNNIK Y, ABLAVATSKI A, GRISHCHENKO I, et al. Real-time facial surface geometry from monocular video on mobile gpus[J]. arXiv:1907.06724, 2019.
- [2] NYISZTOR K. Introduction to Augmented Reality with ARKit[EB/OL]. 2019. <https://www.pluralsight.com/resources/blog/guides/introduction-to-augmented-reality-with-arkit>.
- [3] GINOSAR S, BAR A, KOHAVI G, et al. Learning Individual Styles of Conversational Gesture[C/OL]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 2019: 3497-3506. https://openaccess.thecvf.com/content_CVPR_2019/papers/Ginosar_Learning_Individual_Styles_of_Conversational_Gesture_CVPR_2019_paper.pdf.
- [4] KUCHERENKO T, HABERNAL I, BESKOW J, et al. Multimodal Analysis of the Predictability of Hand-Gesture Properties[J/OL]. Frontiers in Computer Science, 2021, 3: 10. https://web.cs.ucdavis.edu/~neff/papers/kucherenko2021hand_property_predictability_final.pdf.
- [5] LIU H, ZHU Z, IWAMOTO N, et al. BEAT: A Large-Scale Semantic and Emotional Multi-Modal Dataset for Conversational Gestures Synthesis[J]. arXiv preprint arXiv:2203.05297, 2022.
- [6] CASSELL J, VILHJÁLMSSON H H, BICKMORE T. BEAT: the Behavior Expression Animation Toolkit[C/OL]//SIGGRAPH '01: Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques. 2001: 477-486. <https://doi.org/10.1145/383259.383315>. DOI: 10.1145/383259.383315.
- [7] HUANG C M, MUTLU B. Robot behavior toolkit: generating effective social behaviors for robots[C/OL]//HRI '12: Proceedings of the Seventh Annual ACM/IEEE International Conference on Human-Robot Interaction. 2012: 25-32. <https://doi.org/10.1145/2157689.2157694>. DOI: 10.1145/2157689.2157694.
- [8] KIPP M. Gesture generation by imitation: from human behavior to computer character animation[C/OL]//. 2005. <https://api.semanticscholar.org/CorpusID:26271318>.
- [9] WAGNER P, MALISZ Z, KOPP S. Gesture and speech in interaction: An overview

- [J/OL]. Speech Communication, 2014, 57: 209-232. DOI: 10.1016/j.specom.2013.09.008.
- [10] ZHOU Y, BARNES C, LU J, et al. On the Continuity of Rotation Representations in Neural Networks[J]. arXiv preprint arXiv:1812.07035v4, 2020.
- [11] LIU H, ZHU Z, BECHERINI G, et al. Emage: Towards unified holistic co-speech gesture generation via expressive masked audio gesture modeling[J]. arXiv:2401.00374, 2024.
- [12] AO T, ZHANG Z, LIU L. GestureDiffuCLIP: Gesture Diffusion Model with CLIP Latents[J/OL]. ACM Trans. Graph., DOI: 10.1145/3592097.
- [13] CHHATRE K, DANĚČEK R, ATHANASIOU N, et al. AMUSE: Emotional Speech-driven 3D Body Animation via Disentangled Latent Diffusion[C/OL]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 2024: 1942-1953. <https://amuse.is.tue.mpg.de>.
- [14] CHEN J, LIU Y, WANG J, et al. DiffSHEG: A Diffusion-Based Approach for Real-Time Speech-driven Holistic 3D Expression and Gesture Generation[C]//CVPR. 2024.
- [15] ZHANG Z, AO T, ZHANG Y, et al. Semantic Gesticulator: Semantics-Aware Co-Speech Gesture Synthesis[J]. ACM Transactions on Graphics (TOG), 2024, 43(4): 1-17.
- [16] ZHU L, LIU X, LIU X, et al. Taming Diffusion Models for Audio-Driven Co-Speech Gesture Generation[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023: 10544-10553.
- [17] YANG S, WU Z, LI M, et al. DiffuseStyleGesture: Stylized Audio-Driven Co-Speech Gesture Generation with Diffusion Models[C/OL]//Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23. 2023: 5860-5868. <https://doi.org/10.24963/ijcai.2023/650>. DOI: 10.24963/ijcai.2023/650.
- [18] HOGUE S, ZHANG C, DARUGER H, et al. DiffTED: One-shot Audio-driven TED Talk Video Generation with Diffusion-based Co-speech Gestures[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops. 2024: 1922-1931.
- [19] DEICHLER A, MEHTA S, ALEXANDERSON S, et al. Diffusion-Based Co-Speech

- Gesture Generation Using Joint Text and Audio Representation[C/OL]//ICMI '23: Proceedings of the 25th International Conference on Multimodal Interaction. 2023: 755-762. <https://doi.org/10.1145/3577190.3616117>. DOI: 10.1145/3577190.3616117.
- [20] ALEXANDERSON S, KUCHERENKO T, HENTER G E, et al. DiffGesture: Diffusion-based Co-Speech Gesture Generation[C/OL]//Proceedings of the 23rd ACM International Conference on Intelligent Virtual Agents (IVA 2023). 2023: 1-8. <https://dl.acm.org/doi/10.1145/3570945.3607282>. DOI: 10.1145/3570945.3607282.
- [21] YOON Y, KIM S, LEE J, et al. Speech Gesture Generation from the Trimodal Context of Text, Audio, and Speaker Identity[C/OL]//Proceedings of the 2020 ACM International Conference on Multimodal Interaction (ICMI). 2020: 22-30. <https://dl.acm.org/doi/10.1145/3382507.3418838>. DOI: 10.1145/3382507.3418838.
- [22] KENDON A. Gesture: Visible Action as Utterance[M]. Cambridge, UK: Cambridge University Press, 2004.
- [23] MCNEILL D. Hand and Mind: What Gestures Reveal about Thought[M]. Chicago, IL: University of Chicago Press, 1992.
- [24] BAARS S, ANDEWEG B. 'Wapper wat meer met je handen' [J/OL]. Tijdschrift voor Taalbeheersing, 2019, 41(1): 3-17. <https://www.aup-online.com/content/journals/10.5117/TVT2019.1.001.BAAR>. DOI: <https://doi.org/10.5117/TVT2019.1.001.BAAR>.
- [25] HADAR U, BUTTERWORTH B. Iconic gestures, imagery, and word retrieval in speech[J]. Semiotica, 1989, 75(1/2): 63-83.
- [26] ESTEVE-GIBERT N, PRIETO P, PONS X, et al. The timing of head movements: The role of prosodic heads and edges[J/OL]. The Journal of the Acoustical Society of America, 2017, 141(6): 4727-4739. DOI: 10.1121/1.4986649.
- [27] MAHMOOD N, GHORBANI N, TROJE N F, et al. AMASS: Archive of Motion Capture as Surface Shapes[C]//Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). 2019: 5442-5451.
- [28] LOPER M, MAHMOOD N, ROMERO J, et al. SMPL: A Skinned Multi-Person Linear Model[J]. ACM Transactions on Graphics (TOG), 2015, 34(6): 248:1-248:16.
- [29] GLEICHER M. Retargetting Motion to New Characters[C]//Proceedings of the 25th Annual Conference on Computer Graphics & Interactive Techniques (SIGGRAPH '

- 98). ACM, 1998: 33-42.
- [30] MARTINELLI G, GARAU N, BISAGNO N, et al. Skeleton-Aware Motion Retargeting Using Masked Pose Modeling[C]//European Conference on Computer Vision (ECCV) Workshops, LNCS 15624. Springer, 2024: 287-303.
- [31] FU C, WANG Y, ZHANG J, et al. MambaGesture: Enhancing Co-Speech Gesture Generation with Mamba and Disentangled Multi-Modality Fusion[EB/OL]. 2024. <https://arxiv.org/abs/2407.19976>. arXiv: 2407.19976 [cs.HC].
- [32] EKMAN P, FRIESEN W V. Facial Action Coding System: A Technique for the Measurement of Facial Movement[M]. Palo Alto, California: Consulting Psychologists Press, 1978.
- [33] ARKit in iOS - Tracking and Visualizing Faces. 2024. https://developer.apple.com/documentation/arkit/arkit_in_ios/content_anchors/tracking_and_visualizing_faces.
- [34] OZEL M. ARKit to FACS Cheat Sheet[EB/OL]. 2022. <https://melindaozel.com/arkit-to-facs-cheat-sheet/>.
- [35] ALEXANDERSON S, HENTER G E, KUCHERENKO T, et al. Style-Controllable Speech-Driven Gesture Synthesis Using Normalising Flows[C/OL]//Computer Graphics Forum (Proc. Eurographics). 2020: 487-496. <https://people.kth.se/~ghe/pubs/pdf/alexanderson2020style.pdf>.
- [36] KUCHERENKO T, BESKOW J, KJELLSTRÖM H, et al. Moving Fast and Slow: Analysis of Representations and Post-processing in Speech-Driven Gesture Generation[J/OL]. arXiv preprint arXiv:2107.12305, 2021. <https://people.kth.se/~ghe/pubs/pdf/kucherenko2021moving.pdf>.
- [37] BOTTOU L. Large-Scale Machine Learning with Stochastic Gradient Descent[C/OL] //International Conference on Computational Statistics. 2010. <https://api.semanticscholar.org/CorpusID:115963355>.
- [38] Unity Editor [Computer Software][A/OL]. Unity Technologies. <https://unity.com/>.
- [39] Mechanim Animation System [Computer Software Module][A/OL]. Unity Technologies. <https://docs.unity3d.com/Manual/AnimationOverview.html>.
- [40] PICO 4 All-in-One VR Headset [Hardware Device][A/OL]. ByteDance Inc. <https://www.pico-interactive.com/>.
- [41] SCHWIND V, RESCH S, SEHRT J. The HCI User Studies Toolkit: Supporting

Study Designing and Planning for Undergraduates and Novice Researchers in Human-Computer Interaction[C/OL]//Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems (CHI EA '23). ACM, 2023: 7. DOI: 10.1145/3544549.3585890.

- [42] YOON Y, CHA B, LEE J H, et al. Speech gesture generation from the TRIMODAL context of text, audio, and speaker identity[J]. arXiv:2009.02119, 2020.
- [43] LI R, YANG S, ROSS D A, et al. AI Choreographer: Music Conditioned 3D Dance Generation With AIST++[C]//Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). 2021: 13401-13412.

附录 A 手势生成对比视频

为直观展示不同模型在语音驱动手势生成任务中的表现，本文提供了基于 BEAT 数据集的可视化对比结果。具体而言，在说话人 2、4、6、8 的测试语音上，分别对本文模型 FaceCapGes、以及 CaMN^[5] 与 DiffSHEG^[14] 生成的手势序列进行了对比展示。

对应的生成结果视频可通过以下链接访问：

- **Gesture Generation Comparison Videos:**

[https://drive.google.com/drive/folders/1TmPqF-e6eHwSsMCQ2JNyPgGdMzKII_DbW?
usp=sharing](https://drive.google.com/drive/folders/1TmPqF-e6eHwSsMCQ2JNyPgGdMzKII_DbW?usp=sharing)

附录 B 代码与实现资源

本文模型 FaceCapGes 的完整训练与推理代码已开源，以便于复现本文中的实验结果与评估指标。代码仓库地址如下：

- **FaceCapGes GitHub Repository:**

<https://github.com/IORGestureTeam/FaceCapGes>

学术论文和科研成果目录

学术论文

- [1] First Author. FaceCapGes: Real-Time Frame-by-Frame Gesture Generation from Audio, Facial Capture, and Head Pose[C]. Computer Graphics International (CGI 2025), 2025.