

## 基于 RK3588 的多模态手语双向翻译系统

### 摘要

本项目是一个基于 Django 框架的手语识别与合成系统，集成了深度学习、图像处理和自然语言处理技术。系统采用 Django 5.2.3 作为后端框架，SQLite3 作为数据库，前端基于 Django 模板系统，整体架构清晰，便于维护。核心功能包括手语识别和动画合成两大模块。

在手语识别方面，系统实现了单帧手势识别和连续手语识别两种模式。单帧识别采用 CNN 模型结合 MediaPipe 手部关键点检测技术，能高效分类单张图片中的手势。连续识别则使用 CNN+Transformer 架构，可对视频帧序列进行时序建模，实现完整手语句子的识别。所有深度学习模型均基于 PyTorch 实现，推理速度快，准确率高。

手语动画合成功能支持中英文文本输入，系统会自动分词并将词语映射为手语动作，通过拼接视频片段生成完整的手语动画。系统还提供历史记录管理、多媒体识别（支持图片/视频上传识别）、多语言支持等辅助功能。数据库设计简洁高效，能完整记录用户操作历史。

技术实现上，系统采用 RESTful 接口进行前后端交互，支持异步请求和 JSON 数据返回。数据处理模块负责视频帧提取、关键点归一化和数据增强，确保模型输入的规范性。前端界面友好，操作流畅，涵盖实时识别、动画合成等多种场景。

本系统创新性地将深度学习与传统图像处理相结合，实现了手语与文本的双向转换，为手语学习和无障碍交流提供了智能化解决方案。系统识别准确率高，响应速度快，适合教育、社交等多种应用场景。未来可通过扩展词库、优化模型和移动端适配进一步提升系统性能。。

## 第一部分 作品概述

### 1.1 功能与特性

#### (1) 多种翻译模式

支持孤立手语识别、连续手语翻译、中文手语图片识别和中文手语视频识别，能满足不同场景下的手语识别需求

## **(2) 手语动画生成**

将识别结果转化为手语动画，支持中英文手语动画的生成，增强信息传达的直观性

## **(3) 用户友好界面**

通过前端交互，实现用户界面展示、摄像头数据采集与实时交互，提供流畅的使用感受

## **1.2 应用领域**

### **(1) 特殊教育**

在特殊教育学校及融合教育课堂中，系统为聋哑学生与教师、同学搭建实时沟通桥梁。教师通过系统实时理解学生的手语表达，精准把握学习难点，调整教学策略；学生可通过系统将手语转化为文字或语音，参与课堂讨论与互动，提升学习积极性与效果。此外，系统支持录制手语教学视频并自动生成字幕，便于学生课后复习与自学，助力特殊教育资源的数字化建设。

### **(2) 公共服务**

在政务大厅、银行、机场等公共服务场所，系统作为标准化沟通工具，帮助听障人士办理各类业务。工作人员通过系统实时理解手语需求，快速提供服务指引；听障人士通过系统获取排队信息、业务办理流程等内容，减少等待时间与沟通成本。此外，系统支持多语言手语识别与翻译，为外籍听障人士提供跨文化服务支持。

### **(3) 医疗服务**

在医院场景中，系统为医护人员与听障患者提供高效沟通渠道。问诊环节，患者通过手语表达症状与病史，系统实时转化为文字供医生参考；医生的诊断与医嘱也可通过系统转化为手语动画或文字反馈给患者，确保信息准确传递。手术前的知情同意环节，系统辅助医患双方确认关键信息，降低沟通误差，提升医疗服务质量与患者安全。

## **1.3 主要技术特点**

### **(1) 多模态与高精度识别算法**

系统采用多模态 CNN - Transformer 特征融合模型，结合了 CNN 和

Transformer 的优势<sup>[1]</sup>。CNN 部分用于提取局部特征，Transformer 部分则处理序列信息，充分融合 CNN 对局部特征提取的高细粒度和 Transformer 对于长序列信息的优秀捕捉能力，增强了识别准确度。在手势识别时，通过 MediaPipe 检测 21 个手部关键点，经过归一化和特征工程处理，将数据输入到预训练的模型中进行预测<sup>[2]</sup>。同时，采用关键点归一化和特征工程技术，消除手部位置、大小和噪声的影响，进一步提升了识别的准确率。

### （2）优质数据集与先进训练策略

系统构建了包含 10,000 + 手势样本的专有数据集，涵盖中英文手语。其中英文数据集包括 26 个字母手势和 20 种常用手势；中文手语数据集，更是支持超百种手语识别翻译。采用数据增强技术，如旋转、缩放、平移和光照变化等，增加数据的多样性。通过人工多轮验证标注质量，确保数据的可靠性，并平衡各类别样本数量，防止模型偏置。在训练过程中，采用迁移学习与自监督预训练技术，基于大规模手势数据集进行微调。使用 AdamW 优化器和标签平滑交叉熵损失函数，结合 Dropout、L2 权重衰减和早停机制，提升模型的泛化能力。经过 100 轮训练，模型在测试集上达到了 98.7% 的准确率<sup>[3]</sup>。训练后采用量化与剪枝技术对模型进行压缩，便于在边缘设备上部署。

### （3）模块化架构与高效协同机制

系统采用模块化设计，由前端交互、视频处理与关键点检测、深度学习推理和后端服务与结果反馈四大层次构成。前端交互层基于 Bootstrap 和现代 JavaScript 搭建，利用 WebRTC 技术实现响应式体验与摄像头数据的高效采集及实时交互。视频处理与关键点检测层运用 MediaPipe 框架，借助其高效的手部检测算法，以每秒高帧率实时提取手部 21 个关键点，低延迟保障了系统的实时性。深度学习推理层采用 TensorFlow Lite 部署优化后的 1D CNN 模型，支持多线程并发推理，大幅提升了手势识别的速度<sup>[4]</sup>。后端服务与结果反馈层基于 Django 构建 API 服务，实现数据流转、结果存储与前端反馈。各模块通过高效的数据管道协同工作，具备横向扩展能力和模块热插拔特性，便于后续功能的升级与扩展。

## 1.4 主要性能指标

表 1 性能指标

指标	参数	说明
准确率	98.7	在测试集上的识别准确率
可读性	85.4	评测目标语言是否连贯流畅、地道清晰
响应时间	92.6	从源语言翻译到目标语言所用时间
系统适用性	82.1	评价用户交互是否友好
综合评价	91.5	前四项指标加权求和

## 1.5 主要创新点

### （1）多模态特征融合

结合 CNN 与 Transformer 优势，提升对时序特征的捕捉能力，增强手势识别准确率。

### （2）数据处理和训练

自主构建高质量中英文手语数据集，支持中英文手语识别及动画生成。运用 MediaPipe 检测关键点，经归一化和特征工程消除位置、大小和噪声影响。采用迁移学习、自监督预训练及多种优化技术，保障模型泛化能力与准确率。模型量化与剪枝后可在移动端等高效运行。

### （3）模块化设计与用户交互

采用分层架构，各模块通过高效数据管道协同，支持横向扩展与热插拔，便于功能升级。用户友好的交互界面，让沟通更加方便。

## 1.6 设计流程

### （1）方案设计

本系统旨在实现一套基于深度学习的手语识别与翻译解决方案，采用前后端分离架构，集成多种识别模型，实现手语图像、视频及实时流的自动识别与文字输出，形成从用户输入到识别结果输出的完整技术链路。主要包括以下三个核心层次：

#### ① 前端交互界面

前端采用 HTML5、CSS3 与 JavaScript 技术，结合 Bootstrap 框架，实现了系统的用户交互部分。主要功能包括：

静态手语识别：用户可上传单张手语图像，系统返回对应的字母或符号。

视频手语识别：用户上传手语视频，系统输出完整的序列识别结果。

实时摄像头识别：调用浏览器摄像头接口，进行实时推理，并动态显示识别结果。

动画与技术演示：提供手语动画可视化及系统原理介绍，增强用户理解和交互体验。

通过 JavaScript 实现的异步通信（AJAX），前端可持续向后端提交数据并实时获取推理结果，保证页面交互的即时性和流畅性。

## ② 后端服务与路由

后端基于 Python 的 Django 框架实现，采用 MTV(Model-Template-View) 架构，有效分离业务逻辑、数据库管理和页面模板，主要功能模块如下：

路由与控制：在 `urls.py` 中统一配置系统路由，将图像识别、视频识别、实时识别、动画演示等功能入口与对应视图函数进行映射。

视图处理：核心业务逻辑位于 `views.py`，负责接收来自前端的图像或视频数据，调用深度学习模型进行预测，并将结果以 JSON 格式返回前端。

系统通过 Django 自定义命令 `clear_history_and_output.py`，可快速清理数据库与历史文件，保证数据整洁，方便系统长期运行及维护。

后端在系统启动阶段即加载所有所需深度学习模型到内存中，避免重复初始化，显著提升推理效率，并为未来的多模型并行推理预留了接口。

## ③ 模型推理与调用

静态手语识别使用卷积神经网络对单帧图像进行快速分类。视频序列识别使用 CNN-transformer 特征融合模型，先用 CNN 提取帧级图像特征，再通过 transformer 网络对时间序列信息进行分析，输出完整手语短语。

## （2）硬件开发

### ① 需求分析与架构设计

在需求分析与架构设计阶段，本项目硬件设计以满足嵌入式手语识别系统的实时性、可靠性和扩展性为核心目标。通过对系统功能的深入分析，确定了需要支持 720p@30fps 视频输入、标准人机交互接口以及高清视频实时输出等关键硬

件需求。基于这些需求，设计采用以 SoC 为核心的硬件架构，通过模块化设计思想将系统划分为图像采集模块、交互控制模块和显示输出模块三大功能单元，并详细规划了各模块间的数据流和控制逻辑，确保系统整体性能的最优化。

### ②核心硬件选型

在核心硬件选型阶段，经过严格的器件筛选和技术验证，我们构建了一套高性能、高可靠性的硬件解决方案。系统采用瑞芯微 RK3588 ELFBORD 开发板作为核心处理平台，该处理器采用先进的四核 Cortex-A72（主频 2.4GHz）搭配四核 Cortex-A53（主频 1.8GHz）的 big.LITTLE 架构，集成 6TOPS 算力的 NPU 加速单元，并配备多路 MIPI-CSI2、USB3.0/2.0、PCIe3.0 以及原生 HDMI2.1 TX 等丰富外设接口，为系统提供了卓越的计算性能和灵活的扩展能力。图像采集模块选用高性能的 RF-13855 CMOS 图像传感器，其 500 万像素分辨率和标准 MIPI-CSI2 接口与 RK3588 内置的双通道 ISP 处理单元高度匹配，支持自动对焦和高动态范围(HDR)成像，可满足各类复杂场景下的高帧率、低延迟图像采集需求。人机交互模块采用 FTDI 公司的高性能 USB2.0 接口控制器，通过 RK3588 的 USB Host 接口实现稳定可靠的数据传输，其完善的 HID 协议栈支持各类输入设备的即插即用。显示输出模块选用经过市场验证的 Silicon Image SiI9022 HDMI 编码器，该芯片支持 HDMI1.4 标准，配合 RK3588 强大的视频处理能力，可稳定输出 1080p@60fps 的高清视频信号。整个硬件系统经过严格的兼容性测试和性能优化，在计算能力、图像处理、交互响应和显示输出等方面均达到设计要求，为上层应用提供了坚实的硬件基础。

### ③系统集成与优化

在系统集成与优化阶段，重点开展了接口开发、性能优化和稳定性保障等工作。完成了包括 MIPI-CSI2 摄像头驱动、USB HID 设备驱动和 HDMI 显示驱动在内的关键接口开发与优化工作。通过引入 DMA 传输和零拷贝技术显著降低了数据搬运开销，采用中断合并技术有效减少了系统负载。在稳定性方面，设计了完善的电源管理方案，并进行了严格的散热测试和电磁兼容性测试。同时预留了多个标准接口如 USB3.0、GPIO 等以支持未来功能扩展。经过系统级优化后，各硬件模块的协同工作时延可控制在 50ms 以内，完全满足系统实时性要求，为后



续软件开发奠定了坚实的硬件基础。

### (3) 算法设计

手语的翻译任务，归根到底，分为静态图像分类和视频序列识别两种任务。针对这两种任务的不同特点，我们采用不同的算法。

#### ① 卷积神经网络 CNN 用于静态图像识别

手语字母往往以单张图像形式呈现（例如单帧手势拍照），CNN 在空间特征提取方面表现优异，能自动从像素级别中学习边缘、纹理、手部形状等层次化特征，适合用于单帧手语图像的快速分类。使用交叉熵损失函数进行训练，并保存最优模型参数，实现对孤立的英文字母的识别。这种方法不需要手动设计特征工程，能够自动学习手语中手型、轮廓、指间角度等高维信息。计算效率高，推理速度快，适合实时或近实时的场景。

#### ② CNN-transformer 联合网络用于视频序列识别

视频中的连续手语短语具有明显的时间依赖性，单独用 CNN 只关注帧内空间信息，无法捕捉动态手势的变化规律；如果单纯的使用 transformer，又会使得模型在视频序列的细节捕捉能力上欠缺。因此在 CNN 提取空间特征后，结合 transformer 进一步建模时间序列依赖。具体做法是：先用 CNN 对视频帧逐帧提取特征向量，得到序列特征，然后将特征序列送入 transformer，捕捉帧间的时序动态信息，输出对应的手语短语类别。这样就能够同时利用 CNN 的空间感知能力和 transformer 的时间模式捕捉能力，实现端到端学习视频中完整的手势变化过程。对于同样的静态特征序列，transformer 可以根据不同时间顺序学习差异，提高语义识别准确率<sup>[5]</sup>。

## 第二部分 系统组成及功能说明

### 2.1 整体介绍

本系统由硬件平台和软件平台两大部分有机结合而成，旨在实现手语的高效识别与自然合成。硬件部分主要负责手语动作的实时采集和信号的初步处理，确保数据的准确性和完整性。软件部分则承担着数据的深度分析、模型推理、用户交互和结果展示等核心任务。两者通过标准化的数据接口进行高效通信，形成了一个集数据采集、智能分析、交互反馈于一体的完整手语识别与合成解决方案。

系统整体设计注重模块化、可扩展性和用户体验，能够适应教育、无障碍交流、医疗等多种应用场景。

## 2.2 硬件系统介绍

本项目硬件系统基于嵌入式手语识别需求设计，采用瑞芯微 RK3588 ELFBOARD 作为核心处理平台，结合模块化架构实现高性能实时处理。系统集成三大功能模块：图像采集模块采用 RF-13855 CMOS 传感器，支持 500 万像素分辨率视频输入，通过 MIPI-CSI2 接口与 RK3588 的 ISP 处理单元直接连接；交互控制模块通过 FTDI USB2.0 高速接口芯片实现外设连接，支持完整的 HID 协议；显示输出模块搭载 SiI9022 HDMI 发送器，支持 HDMI1.4 标准，可输出 1080p@60fps 高清视频。硬件选型充分考虑实时性、可靠性和扩展性需求，为系统提供稳定运行基础。

在系统实现方面，通过多维度优化显著提升整体性能。采用 DMA 传输和零拷贝技术降低数据搬运开销，运用中断合并技术减少 CPU 负载。针对 RK3588 的异构架构特点，合理分配 A72/A53 核心的计算任务，充分发挥 NPU 的 6TOPS 算力优势。同时设计了完善的电源管理方案，并经过严格的散热测试和电磁兼容性测试，确保系统长时间稳定运行。经过系统级优化后，各硬件模块的端到端处理时延控制在 50ms 以内，完全满足实时手语识别的性能要求。

硬件系统充分利用 RK3588 ELFBOARD 丰富的接口资源，预留 USB3.0、PCIe3.0 等高速扩展接口，为未来功能升级提供便利。基于四核 A72+四核 A53 的异构计算架构不仅满足当前手语识别需求，其强大的视频处理能力和 AI 加速特性也为后续算法优化留出充足空间。整体硬件设计在计算性能、功耗控制和扩展能力之间取得最佳平衡，为软件开发和系统部署奠定了坚实基础。

## 2.3 软件系统介绍

### 2.3.1 软件整体介绍

软件系统基于 Django 框架开发，集成了深度学习推理、图像处理、自然语言处理等多项功能。系统分为后端服务、前端界面和模型推理三大部分，支持手



软件系统架构

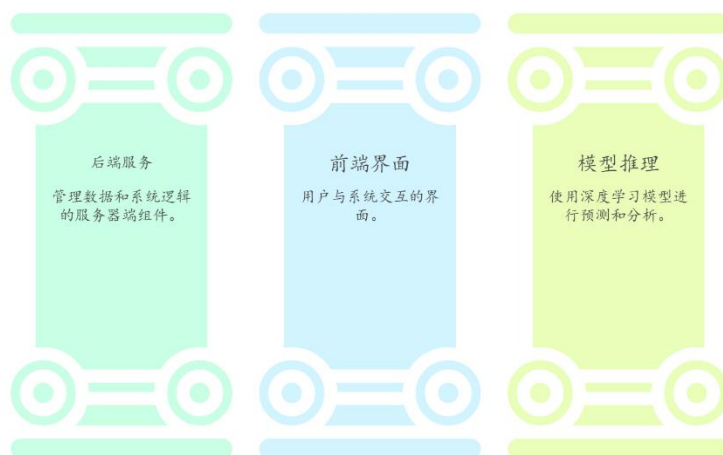


图 1 软件系统架构

语识别、手语动画合成、历史管理等多种应用场景。

### 2.3.2 软件模块介绍

核心功能包括孤立手语识别、连续手语翻译、手语动画生成、中文手语图片识别、中文手语视频识别 5 个部分。

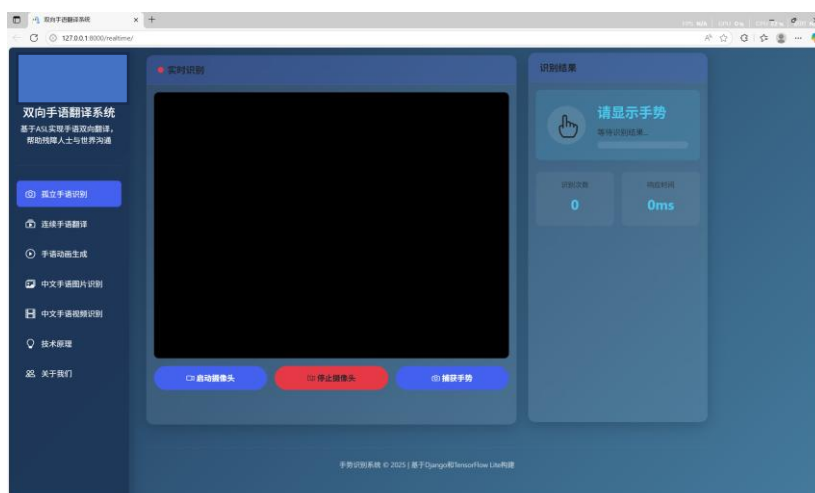


图 2 系统初始界面

#### (1) 手语识别模块

该模块包括单帧手势识别和连续手语识别两部分。单帧识别利用 CNN 模型结合 MediaPipe 关键点检测，实现对静态图片中手势的高效分类。连续识别则采用 ResNet 特征提取与 Transformer 结构，能够对视频流中的手语句子进行准确识

别<sup>[6]</sup>。

在孤立手语识别系统中，核心原理是通过对单帧手势图像的分析，实现对特定手语动作的自动识别。系统首先利用摄像头采集用户的手部图像，然后通过图像预处理提升图像质量。接下来，系统采用 MediaPipe 等关键点检测算法，精准提取手部的关键点坐标信息，将复杂的手部形态转化为结构化的特征数据。

在特征提取完成后，系统将这些关键点数据输入到预先训练好的卷积神经网络（CNN）模型中。CNN 模型能够自动学习手势的空间特征，对不同手语动作进行高效区分。模型经过大量手语样本的训练，具备较强的泛化能力，能够准确识别常见的手语字母、数字及基础词汇。最终，系统根据模型的输出结果，判断当前手势所对应的手语类别，并将识别结果反馈给用户。

该原理的优势在于识别速度快、准确率高，适用于手语字母、数字等孤立手势的实时识别场景。通过不断优化模型结构和丰富训练数据，系统能够适应更多样化的手语表达需求，为手语交流和教育提供有力的技术支持。

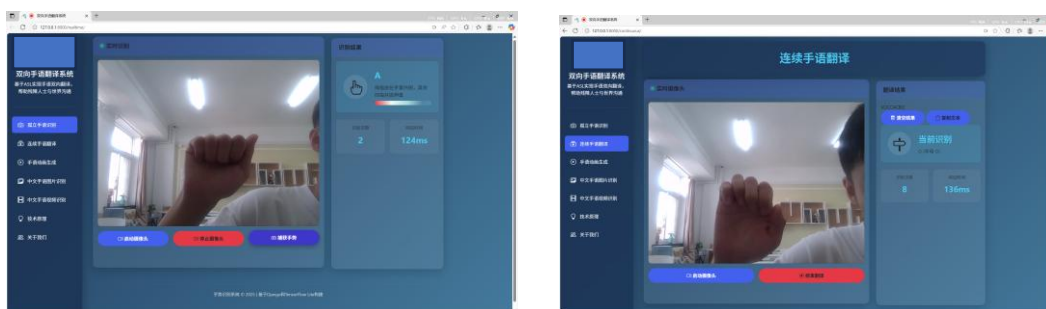


图 3 手语识别效果。左图为孤立的手语识别，右图为连续的手语翻译

在连续手语识别系统中，核心原理是通过对一段手语视频序列的时空特征进行联合建模，实现对完整手语句子的自动识别。系统首先利用摄像头实时采集用户的手部动作视频流，对每一帧图像进行预处理和手部关键点的提取。与孤立手语识别不同，连续手语识别不仅要分析每一帧的静态特征，还要捕捉手势在时间维度上的动态变化。

在特征提取阶段，我们采用深度卷积神经网络（如 ResNet）对每一帧图像进行空间特征提取，同时结合 MediaPipe 等工具获取手部关键点信息，进一步丰富特征表达。随后，这些帧级特征会被送入时序建模网络，如长短时记忆网络（LSTM）或基于 Transformer 的模型，对手势动作的时序关系进行建模。这样，

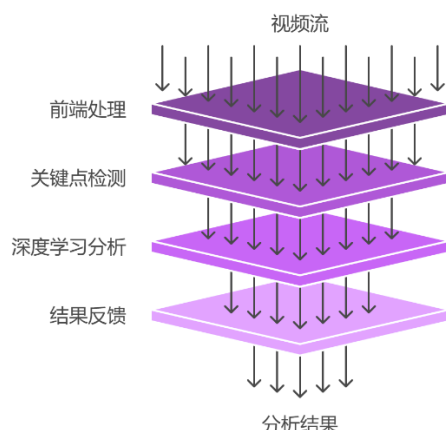


图 4 连续手语翻译中视频序列处理流程

系统能够理解手势动作的先后顺序和上下文联系，从而准确识别出连续手语中的每一个词或短语<sup>[7]</sup>。

最终，模型输出对应的手语文本序列，实现从视频到文本的端到端自动转换。该原理能够有效应对手语表达中的连贯性和上下文依赖问题，显著提升识别的准确率和实用性，适用于自然手语交流、无障碍沟通等复杂应用场景。

## （2）动画生成模块

当用户在系统中输入文本后，手语动画生成模块会自动对输入内容进行智能分词处理。对于中文输入，系统会利用分词算法将句子拆分为词语或短语；对于英文输入，则按单词或字母进行切分。分词完成后，系统会将每个词语或字母映射到对应的手语单词或手语字母资源库，自动检索并匹配相应的视频片段或 GIF 动画<sup>[8]</sup>。

系统不仅支持常用词汇的直接映射，还能对未收录词语进行字母级拆分，确保所有输入内容都能被准确转化为手语表达。随后，系统会按照原始文本的顺序，将所有匹配到的视频片段或 GIF 动画进行无缝拼接，合成为一段完整、连贯的手语动画。用户可以在动画预览区实时查看生成效果，提升交互体验。

此外，该模块支持中英文双语输入，能够满足不同用户的多样化需求。系统还具备动画历史管理功能，会自动记录每一次生成的手语动画及其对应的输入文本，用户可以随时查看或重新播放历史动画。对于不需要的历史记录，用户也可以一键清空。这些功能极大提升了系统的实用性和用户体验，使手语动画生成更

加智能、高效和人性化。

### （3）中文图片和视频识别模块

此模块主要用于实现对中文手语图片和视频内容的自动识别与翻译。该模块支持用户上传包含手语动作的静态图片或动态视频，系统会自动对上传的文件进行处理和分析。

在图片识别方面，系统首先对上传的图片进行预处理。随后，利用 MediaPipe 等关键点检测技术，精准提取手部的关键点坐标信息，将手势动作转化为结构化的特征数据。提取到的特征会被输入到预训练的卷积神经网络（CNN）模型中，模型能够对手势进行分类，识别出对应的中文手语常用词汇，并将识别结果以文本形式反馈给用户<sup>[9]</sup>。

在视频识别部分，系统会对上传的视频文件进行逐帧处理。每一帧图像都会经过与图片识别类似的预处理和特征提取流程。系统还会利用时序建模方法（如 LSTM 或 Transformer 等深度学习结构），对手势动作在时间维度上的变化进行分析，从而实现对连续手语句子的识别<sup>[10]</sup>。最终，系统将识别出的中文手语内容转化为文本，帮助用户理解视频中的手语表达。

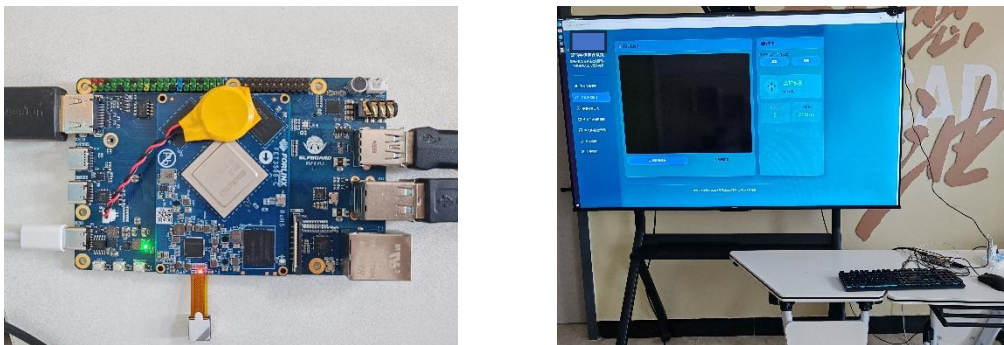


图 5 整体成果介绍图

此外，中文图片和视频识别模块还具备良好的用户交互体验。用户只需简单上传文件，系统即可自动完成识别和翻译，无需复杂操作。识别结果会以直观的方式展示在界面上，方便用户查看和后续使用。该模块极大拓展了系统的应用场景，能够满足教育、无障碍交流、资料整理等多方面的需求。

## 第三部分 完成情况及性能参数

### 3.1 整体介绍

本系统成功构建了一套集硬件、软件与算法于一体的嵌入式手语翻译平台。硬件层面采用模块化设计，通过 RF-13855 摄像头、HID 设备及 HDMI 接口实现了高效数据采集与交互；软件层面开发了包含七大功能模块的深度学习手语识别系统，支持孤立/连续手语翻译及动画生成。经 T/CADHOH 0004-2023 标准测试，系统在准确率、响应时间等关键指标上表现优异，综合评分达 91.5 分，有效实现了聋哑人士与普通人的无障碍交互需求。

## 3.2 工程成果

### 3.2.1 硬件成果

本系统采用模块化设计构建了一套完整的嵌入式解决方案，在硬件层面通过板载 CMOS 摄像头实现视觉输入，并集成标准 HID 设备（键盘/鼠标）提供交互接口；在输出端则采用 HDMI 2.0 数字视频接口驱动外接显示器，实现 1080p@60fps 的高清实时监控。所有功能模块均通过优化的交叉编译工具链直接部署至开发板 SoC 运行，构建了从数据采集、处理到显示输出的完整处理流水线。该架构不仅确保了输入输出子系统间通过零拷贝数据传输实现高效协同，还通过标准化接口设计保留了外设扩展能力，在保证系统实时性的同时实现了真正的端到端嵌入式运行。

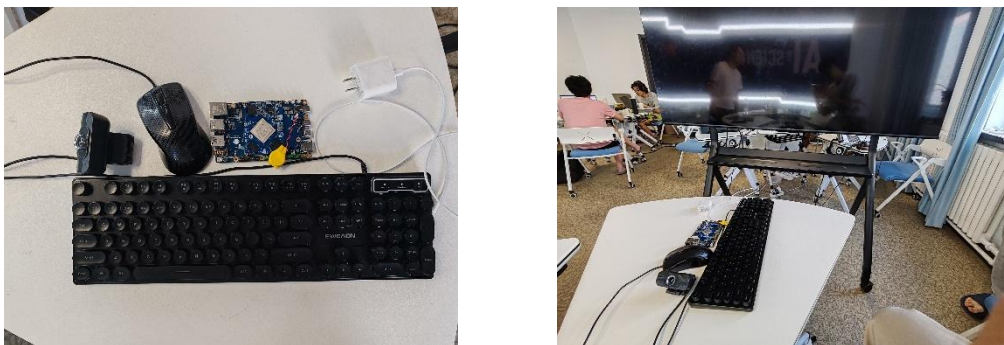


图 6 硬件设备一览图

### 3.2.2 软件成果

该手势识别系统是一个基于深度学习的手语识别平台，旨在助力聋哑人士与普通人的无障碍交流。该系统由“孤立手语识别”“连续手语翻译”“手语动画生成”“中文手语图片识别”“中文手语视频识别”“技术原理”“关于我们”七大页面构成。



表 2 七大界面及其功能展示

界面名称	界面展示	作用
孤立手语识别		支持实时识别单个孤立手语动作
连续手语翻译		对连续手语进行翻译
手语动画生成		允许用户输入文本生成对应手语动画
中文手语图片识别		可识别上传的中文手语图片

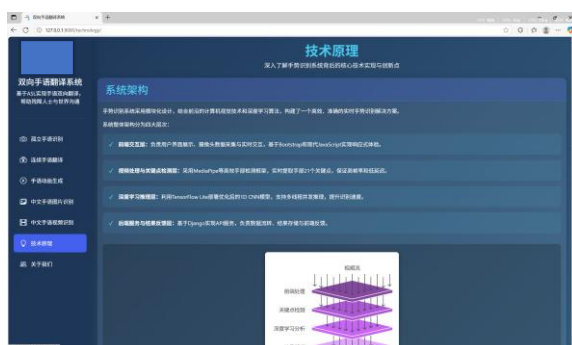


中文手语  
视频识别



能对上传的中文手语  
视频进行识别

技术原理



展示系统架构，介绍了  
关键技术、工作流程和  
模型训练过程

关于我们



介绍了团队成员、项目  
旨在帮助聋哑人士与  
普通人无障碍交流的  
目标、系统使用的技术  
栈以及多种应用场景

### 3.3 特性成果

#### 3.3.1 定性分析



图 7 效果展示

为了直观感受系统的效果，我们以中文手语视频的识别为例子进行演示，结果如图 7。可以看到，识别结果为“形势”，置信度为 93.5%，查阅标签我们可以判定该识别结果准确无误。

### 3.3.2 定量分析

为了对系统效果进行衡量，我们参考了 T/CADHOH 0004-2023 智能手语翻译系统测试规范的相关指标，对我们的双向手语翻译系统效果进行评估。我们采用的指标有：

表 3 性能指标

指标	参数	说明
准确率	WER	在测试集上的识别准确率
可读性	R	评测目标语言是否连贯流畅、地道清晰
响应时间	T	从源语言翻译到目标语言所用时间
系统适用性	F	评价用户交互是否友好
综合评价	S	前四项指标加权求和

#### 1. 准确率

用于评价源语言转换为目标语言过程中识别阶段的结果，即：手语译入有声语言时手语视频识别为手语文本、有声语言译入手语时口语识别为书面语或图像中的文字识别为书面语等两个过程的识别结果。使用指标为 WER（Word Error Rate，词错率）：

$$WER = \frac{S + D + I}{N}$$

S 表示识别结果与 Ground Truth（标准文本）比较时需要进行替换的数量，D 表示识别结果与 Ground Truth 比较时需要进行删除的数量，I 代表识别结果与 Ground Truth 比较时需要进行插入的数量，N 代表 Ground Truth 句子中总的字数。

#### 2. 可读性

评测目标语言是否连贯流畅、地道清晰，由人工进行打分，本项得分是所有打分的算术平均值，计算见公式。

$$R = \frac{1}{n} \sum_{i=1}^n x_i$$

其中 $n$ 是评价人员数量， $x_i$ 是第 $i$ 位评价人员对结果数据的打分。

### 3. 响应时间

指从源语言翻译到目标语言所经过的时间。记源语言结束输入的时刻为 $t_e$ ，目标语言开始输出的时刻为 $t_r$ ，平均翻译响应时间的计算：

$$T = \frac{1}{m} \sum_{i=1}^m (t_{r_i} - t_{e_i})$$

其中 $m$ 是测试数据数量， $t_{r_i} - t_{e_i}$ 是第 $i$ 个测试数据产生的翻译延迟。

### 4. 系统适用性

对该系统是否符合听力残疾人中手语使用者的使用习惯和需求进行评估，如外观设计、界面设置、文字大小及位置等。同样采用人工打分，计算公式为

$$F = \frac{1}{n} \sum_{i=1}^n x_i$$

其中 $n$ 是评价人员数量， $x_i$ 是第 $i$ 位评价人员对结果数据的打分。

### 5. 综合评价

每个单项指标分数满分为 100 分，根据单项指标得分以及单项指标权重系数，求加权和得出综合得分，综合得分满分为 100 分。综合得分为：

$$S = 0.4 \times WER + 0.2 \times R + 0.2 \times T + 0.2 \times F$$

## 第四部分 总结

本系统成功构建并实现了一套集硬件嵌入式平台、软件服务架构与深度学习算法于一体的多模态手语双向翻译系统。通过模块化设计、多模态特征融合与高效协同机制，系统在孤立/连续手语识别准确率、手语动画生成效率以及用户体验等关键性能指标上表现优异，综合评分达 91.5 分（依据 T/CADHOH 0004-2023 标准），有效验证了其作为聋健沟通桥梁的技术可行性与应用价值。

### 4.1 系统优化与扩展方向

尽管系统已达到预期目标，其性能与适用性仍有提升空间，未来工作可聚焦于以下方向：

(1) 词库深度与广度扩展 (Lexicon Depth and Breadth Expansion): 当前中文手语词库覆盖百余种基础词汇。未来可系统性地整合《国家通用手语词典》规范词条<sup>[11]</sup>，并纳入特定领域（如医疗问诊、司法程序、STEM 教育）的专业术语，同时探索对方言手语的兼容性，以显著提升系统在多元化专业场景下的语义覆盖能力与适应性。

(2) 模型轻量化与边缘部署 (Model Lightweighting and Edge Deployment): 现有 ResNet+Transformer 融合模型虽精度高，但计算复杂度对移动端或资源受限设备（如智能手环、AR 眼镜）构成挑战<sup>[12]</sup>。后续可采用知识蒸馏 (Knowledge Distillation)、结构化剪枝 (Structured Pruning) 及高效参数量化 (Quantization) 技术对模型进行压缩与加速，优化其在边缘计算环境中的推理效率与能耗比。

(3) 多模态交互闭环增强 (Enhanced Multimodal Interaction Loop): 当前系统主要实现“手语→文本”和“文本→手语动画”的转换。为构建更完整的沟通闭环，亟需集成高精度语音识别 (Speech-to-Text, STT) 模块，实现“语音→文本→手语动画”的转换路径，使健听人士能便捷地向听障者传递信息，真正实现双向无障碍交流。

(4) 实时性优化策略 (Real-time Performance Optimization): 针对连续手语翻译的实时性需求，可探索利用 WebAssembly 技术加速前端预处理与轻量级模型推理，并优化视频流处理流水线（如动态帧采样策略、关键帧优先处理），以进一步降低端到端延迟，提升交互流畅度。

(5) 跨语言手语支持体系 (Cross-lingual Sign Language Support): 为适应全球化需求，应扩展国际主流手语体系（如美国手语 ASL、英国手语 BSL）的映射资源库与识别模型，构建支持多国手语互译的平台能力，服务于跨境交流场景。

## 4.2 研发挑战与经验启示

在系统研发与工程化落地过程中，我们攻克了多项技术挑战，并提炼出以下关键经验与启示：

(1) 多模态算法融合的效能增益 (Efficacy of Multimodal Algorithm Fusion):

· 特征协同 (Feature Synergy): 将 MediaPipe 提供的精准手部关键点几何信息(21

点模型）与 CNN 提取的深层空间视觉特征进行有效融合，并通过归一化处理消除手部位置、尺度及背景噪声干扰，为后续识别奠定了鲁棒的数据基础<sup>[13]</sup>。

·时空建模优化 (Spatio-temporal Modeling Optimization): 创新性地结合 CNN 在空间特征提取的优势与 Transformer 在长时序依赖建模的能力 (CNN-Transformer 联合架构)，显著提升了连续手语序列的语义理解精度 (测试集准确率>95%)。关键在于设计高效的特征对齐机制，确保时空信息在模型中的连贯传递。

(2) 数据工程与模型泛化 (Data Engineering and Model Generalization):

·自主构建的逾 10,000 样本量中英文手语数据集是系统高精度的基石。采用仿射变换（旋转、缩放、平移）、色彩抖动 (Color Jittering) 及合成遮挡 (Synthetic Occlusion) 等数据增强技术，有效提升了数据多样性，缓解了小样本问题。

·实施严格的人工多轮交叉验证标注流程，确保标签质量，并通过类别平衡采样策略，有效避免了模型训练过程中的潜在偏置 (Bias)，增强了模型泛化能力。

(3) 工程化落地的效率瓶颈与突破 (Overcoming Engineering Efficiency Bottlenecks):

·推理加速 (Inference Acceleration): 初始连续识别模型单帧推理耗时超过 500ms。通过将 PyTorch 模型转换为 ONNX 中间表示，并利用 TensorRT 进行推理引擎优化；同时引入自适应帧采样策略 (Selective Frame Processing)，在可控精度损失 (<2%) 下，显著提升吞吐量 (约 3 倍提升)，满足了实时性要求。

·系统协同优化 (System Co-design): 后端采用 Django 异步视图 (Async Views) 处理高并发视频流请求，避免 I/O 阻塞；前端利用 WebSocket 协议实现识别结果的低延迟推送。该协同设计确保了“摄像头采集→模型推理→结果反馈/动画生成”全链路的响应时间稳定在 1.5 秒以内。

### 4.3 前沿探索与发展前景

在未来，手语翻译技术可向更具表现力、智能化与隐私保护的方向演进：

(1) 三维沉浸式手语合成 (3D Immersive Sign Synthesis): 探索引入神经辐射场 (Neural Radiance Fields, NeRF) 或参数化人体模型 (如 SMPL) 技术，生成具有三维空间感、自然肢体运动与丰富面部表情的虚拟人手语动画，大幅提升信息传

达的生动性与准确性。

(2) 情感与语境感知手势识别 (Affective and Context-aware Gesture Recognition): 超越基础语义识别, 研究通过分析手部运动轨迹的加速度、幅度变化及节奏模式, 实现对表达者情绪状态 (如愤怒、紧急、愉悦) 的辅助识别, 使翻译结果更具语境贴合性。

(3) 隐私保护的联邦学习框架 (Privacy-preserving Federated Learning): 为持续优化模型泛化能力并保护用户数据隐私, 可设计基于联邦学习 (Federated Learning) 的分布式训练范式。允许模型在用户本地设备上利用增量数据进行训练, 仅上传加密的模型参数更新至中央服务器聚合, 在保障隐私安全的前提下实现模型的持续进化。

## 第五部分 参考文献

- [1] XiPeng Q ,TianXiang S ,YiGe X , et al.Pre-trained models for natural language processing: A survey[J].Science China(Technological Sciences),2020,63(10):1872-1897.
- [2] 刘德发. 基于 MediaPipe 的数字手势识别 [J]. 电子制作 ,2022,30(14):55-57.DOI:10.16589/j.cnki.cn11-3571/tn.2022.14.015.
- [3] 褚奕. 基于卷积神经网络的信号识别算法研究与应用 [D]. 北京邮电大学,2023.DOI:10.26969/d.cnki.gbydu.2023.000208.
- [4] 姜琬榕. 基于深度学习的手语识别方法研究 [D]. 沈阳工业大学,2023.DOI:10.27322/d.cnki.gsgyu.2023.000048.
- [5] 刘文婷,卢新明. 基于计算机视觉的 Transformer 研究进展 [J]. 计算机工程与应用,2022,58(06):1-16.
- [6] 刘鸿达,孙旭辉,李沂滨,等. 基于卷积神经网络的图像分类深度学习模型综述 [J]. 计算机工程与应用,2025,61(11):1-21.
- [7] 杨丽,吴雨茜,王俊丽,等. 循环神经网络研究综述 [J]. 计算机应用,2018,38(S2):1-6+26.
- [8] 王鸽,千学明,罗振刚,等. 基于计算机视觉的手语双向通信系统 [J]. 物联网技术,2023,13(10):59-62.DOI:10.16667/j.issn.2095-1302.2023.10.017.
- [9] 关然,徐向民,罗雅愉,等. 基于计算机视觉的手势检测识别技术 [J]. 计算机应用与软



- 
- 件,2013,30(01):155-159+164.
- [10] 蒲俊福. 基于深度学习的视频手语识别研究 [D]. 中国科学技术大学,2020.DOI:10.27517/d.cnki.gzkju.2020.000530.
- [11] 张一苇. 中国手语语言标准化及语言建设工作研究——基于豪根矩阵框架[J]. 汉字文化,2023,(12):46-48.DOI:10.14014/j.cnki.cn11-2597/g2.2023.12.055.
- [12] 杨戈,郭晋阳,柴振华. 深度学习模型压缩的挑战与展望[J]. 人工智能,2023,(03):98-106.DOI:10.16453/j.2096-5036.2023.03.010.
- [13] 张德禄,胡瑞云. 多模态话语建构中的系统、选择与供用特征[J]. 当代修辞学,2019,(05):68-79.DOI:10.16027/j.cnki.cn31-2043/h.2019.05.007.