

HAR 테스트베이스앱장기구축전략: 모듈형실험플랫폼설계

1. 시스템아키텍처설계원칙

1.1 핵심설계철학

1.1.1 모듈러플러그인아키텍처 HAR 테스트베이스앱의근간이되는설계철학은각기능을독립적이고교체가능한플러그인모듈로구성하는것입니다. 이는 “독창적인 HAR 알고리즘을버전별라이브러리로축적”하고 “추후정식서비스에사용”하는핵심목표를달성하기위한필수조건입니다. 포즈추정, 객체추적, 행동분류, 경보시스템등각단계가명확한인터페이스를통해상호작용하도록설계하면, 특정알고리즘의변경이전체시스템에미치는영향을최소화할수있습니다.

구체적인구현방식으로는 Python 의추상기본클래스 (**ABC, Abstract Base Class**) 를활용하여플러그인인터페이스를정의합니다. 예를들어, PoseEstimator 추상클래스는 estimate(frame)-> Skeleton 메서드를강제하며, YOLOv8PoseEstimator, MediaPipePoseEstimator, OpenVINOPoseEstimator 등의구체적구현체가이를상속받습니다. 마찬가지로 Tracker, ActionClassifier, AlertNotifier 등의인터페이스도정의하여, 웹앱의설정패널에서드롭다운메뉴로알고리즘을즉시교체할수있게합니다. 이러한플러그인시스템은 A/B 테스트를용이하게하며, 새로운알고리즘개발시기존인프라를재사용하여개발주기를단축시킵니다.

동적모듈로딩메커니즘도핵심기능으로, 개발자가새로운알고리즘을미리정의된인터페이스를구현하는 Python 클래스로패키징하여특정디렉토리에배치하면, 웹앱이자동으로인식하고실험옵션에추가할수있어야합니다. 이과정에서코드수정이나앱재시작없이새로운알고리즘을통합할수있는핫플러그기능은연구생산성을획기적으로향상시킵니다.

1.1.2 하드웨어추상화계층 (HAL) M1 Pro 에서 NVIDIA Jetson Orin Super, 그리고 Windows PC 까지지원하는다양한하드웨어환경을고려할때, 하드웨어추상화계층 (HAL) 의도입은필수적입니다. HAL 은상위레벨의애플리케이션코드가특정하드웨어세부사항에의존하지않도록하며, 동일한 Python 코드가 Apple Silicon 의 MPS, NVIDIA 의 CUDA/TensorRT, Intel 의 OpenVINO, 순수 CPU 등다양한백엔드에서실행될수있게합니다.

HAL 의핵심구성요소는 InferenceBackend 추상클래스입니다. 이를래스는 initialize(), inference(model, input), get_optimal_batch_size(), memory_stats() 등의메서드를제공하며, 구체적구현체인 MPSBackend, CUDABackend, TensorRTBackend, CPUBackend, OpenVINOBBackend가하드웨어특화된최적화를담당합니다. M1 Pro 환경에서는 PyTorch 의 mps 디바이스를활용하되, YOLOv8-Pose 의경우공식문서에따르면 MPS 에서훈련시알려진버그가있어추론단계에서만 MPS 를사용하고, 안정성이우선될때는 CPU 풀백을제공하는것이바람직합니다. 반면 Jetson 환경에서는 TensorRT 최적화된엔진을로드하여 FP16 또는 INT8 정밀도로추론하는경로가필요합니다.

자동백엔드선택및풀백메커니즘도 HAL 의중요한기능입니다. 시스템초기화시사용가능한하드웨어를자동탐지하여최적의백엔드를선택하고, 지정된백엔드로딩실패시자동으로하위호환가능한백엔드로전환하는기능을제공합니다. 예를들어 Jetson 환경에서 TensorRT 엔진로딩에실패할경우 ONNX Runtime 기반의 CUDA 백엔드로, 그마저실패할경우순수 CPU 백엔드로우아하게성능을저하시키면서도서비스를유지하는것입니다.

1.1.3 알고리즘버전관리및라이브러리화 “독창적인 HAR 알고리즘을버전별라이브러리로축적” 하는것은단순한코드저장을넘어, 체계적인버전관리와재현가능한실험환경을의미합니다. 각알고리즘버전은 **Semantic Versioning(semver)** 을따르며, 메이저. 마이너. 패치형식으로버전을관리합니다. 예를들어, YOLOv8-Pose 통합모듈은 har_pose_yolov8 v1 .2.3과같이명시되며, 각버전은특정모델가중치파일, 전처리파라미터, 후처리로직의조합을고유하게식별합니다.

버전관리의실무적구현으로는, 각알고리즘버전을독립적인 Git 태그로관리하고, 해당버전의메타데이터를 JSON 형식으로기록합니다. 메타데이터에는모델아키텍처, 훈련데이터셋, 입력해상도, FPS 벤치마크결과, 검증정확도등이포함됩니다. 웹앱의 “알고리즘선택” 인터페이스에서는이러한메타데이터를시각적으로표현하여, 사용자가성능-정확도트레이드오프를정보에기반하여결정할수있게합니다. 또한, 동일한입력영상에대해여러버전의알고리즘을병렬실행하여결과를비교하는 “**알고리즘경쟁 (Algorithm Competition)**” 모드를제공하면, 새로운알고리즘의도입결정에객관적인근거를제공할수있습니다.

라이브러리화의최종목표는 pip install har-core와같은명령으로핵심기능을설치하고, from har_core import YOLOv8PoseEstimator, LSTMActionClassifier와같이임포트하여사용할수있는수준입니다. 이를위해 setup.py 또는

`pyproject.toml`을 통한 패키지 메타데이터 정의, PyPI 배포 자동화 파이프라인 구축, 그리고 Sphinx 기반 API 문서화가 필요합니다. 특히 “추후 정식 서비스에 사용하는 것이 목표”라는 점을 고려할 때, 라이선스 명확화 (Apache 2.0 또는 MIT 권장) 와 보안 취약점 스캔을 CI/CD 파이프라인에 통합하는 것이 장기적으로 필수적입니다.

1.1.4 확장 가능한 멀티 카메라 파이프라인 초기 단일 카메라 입력에서 2-4 대 카메라 학장을 염두에 두는 파이프라인 설계는 시스템의 생명력을 결정합니다. 단순히 4 개의 독립적인 싱글 카메라 프로세스를 병렬 실행하는 것이 아니라, 공유 리소스 관리와 부하 분산을 고려한 통합 파이프라인이 필요합니다.

핵심 설계 패턴으로는 **생산자-소비자 모델** 기반의 비동기 프레임 처리를 권장합니다. 각 카메라 입력은 별도의 스레드 또는 프로세스에서 GIL(Global Interpreter Lock)을 우회하는 방식으로 프레임을 획득하여 공유 버퍼 (Queue)에 배치합니다. 중앙 처리 엔진은 이 버퍼로부터 프레임을 소비하여 추론을 수행하고, 결과를 다시 공유 상태 저장소에 기록합니다. Python의 `multiprocessing.Queue`와 `ProcessPoolExecutor`를 활용하거나, 더 고성능이 필요한 경우 ZeroMQ나 Redis를 메시지 브로커로 도입할 수 있습니다.

GPU 메모리 관리는 멀티 카메라 환경에서의 핵심 과제입니다. 4 개의 1080p 스트림을 동시에 처리할 때, 각 스트림마다 별도의 모델 인스턴스를 로드하는 것은 현실적이지 않습니다. 대안으로 **배치 추론 (Batch Inference)**을 활용하여 여러 프레임을 하나의 배치로 묶어 단일 forward pass로 처리하거나, 최소한 모델 가중치를 공유하면서 입력 텐서만 교체하는 방식을 구현해야 합니다. M1 Pro의 16GB 통합 메모리와 Jetson Orin Super의 메모리 제약을 고려할 때, 4 스트림 배치 처리는 충분히 실현 가능하며, 이를 위한 파이프라인 구조는 초기 설계 단계부터 반영되어야 합니다.

설계 원칙	핵심 목표	구현 방식
모듈러 플러그인	알고리즘 교체 용이성	ABC 기반 인터페이스, 동적 로딩
HAL	하드웨어 독립성	Inference Backend 추상화, 자동 선택
버전 관리	실험 재현성	Semantic Versioning, 메타데이터 기록
멀티 카메라	수평 확장성	생산자-소비자 패턴, 배치 추론

1.2 단계별 마이그레이션 전략

1.2.1 Phase 1: M1 Pro 기반 단일 카메라 프로토 타입 첫 번째 단계의 목표는 핵심 알고리즘 파이프라인의 검증과 개발 환경 확립입니다. M1 Pro MacBook 의 자원을 최대한 활용하여 로컬에서 실시간 추론이 가능한 구조를 구축하며, 이 과정에서 HAL의 초기 버전과 기본적인 웹 인터페이스를 개발합니다.

구체적인 하드웨어 활용 전략은 다음과 같습니다. **CPU**는 데이터 전처리 및 웹 서버 구동을 담당하며, Python의 싱글 스레드 성능이 우수한 특성을 활용합니다. **GPU/Neural Engine**은 Metal Performance Shaders(MPS)를 통해 딥러닝 모델 가속을 수행하며, PyTorch 2.1 이상에서 `device='mps'` 설정으로 네이티브 지원됩니다. 입력 장치로는 기본적으로 MacBook Facetime HD Camera를 테스트용으로 사용하고, 확장을 위해 USB 웹캠 (로지텍 C920 등) 또는 RTSP 지원 IP 카메라를 연결하여 실제 CCTV 환경을 모사합니다.

이 단계에서 반드시 해결해야 할 기술적 과제는 **M1 환경에서의 안정적인 실시간 추론**입니다. YOLOv8-Pose의 CPU 모드에서 “Illegal instruction: 4” 오류가 발생할 수 있으며, 이는 MPS 가속을 명시적으로 지정하거나 ONNX Runtime으로의 변환을 통해 우회해야 합니다. 또한 30fps 처리가 불가능한 경우를 대비한 프레임 드롭 로직과, 추론 지연으로 인한 버퍼 오버플로 우방지 메커니즘을 구현해야 합니다. 개발 산출물로는 Streamlit 기반의 기본 대시보드, YOLOv8-Pose를 활용한 실시간 스켈레톤 추출, 그리고 규칙 기반의 간단한 행동 감지 (쓰러짐 등) 가 포함됩니다. 이 단계의 성공 기준은 1080p 입력에서 15fps 이상의 안정적인 처리와 웹 인터페이스를 통한 실시간 모니터링 가능 여부입니다.

1.2.2 Phase 2: 멀티카메라확장 (2-4 대 USB/RTSP) 두번째단계에서는 단일카메라파이프라인을 2-4 대 카메라로 확장하면서, 리소스경쟁과 동기화문제를 해결합니다. USB 웹캠과 RTSP IP 카메라를 혼합하여 사용할 수 있는 유연한 입력 관리 시스템을 구축하며, DeepSORT 기반 다중 객체 추적을 통합합니다.

핵심 기술 과제는 **다중 입력 스트림의 효율적 관리**입니다. USB 웹캠과 RTSP IP 카메라는 기본적으로 다른 I/O 특성을 가지며, 전자는 DirectShow/V4L2를 통해 후자는 FFmpeg/GStreamer를 통해 접근됩니다. HAL의 입력력 장치 모듈은 이러한 차이를 추상화하여 상위레이어에 통일된 프레임 인터페이스를 제공해야 합니다. 특히 RTSP 스트림의 네트워크 지연과 패킷 손실에 대한 견고한 처리, 그리고 USB 대역폭 제한으로 인한 다중 웹캠의 해상도/프레임레이트 조정이 필요합니다.

이 단계에서도 입되는 핵심 알고리즘은 **다중 객체 추적 (Multi-Object Tracking, MOT)**입니다. 단일 프레임 내에서 여러 사람을 검출하는 것을 넘어, 시간에 걸쳐 각 개인의 ID를 일관되게 유지하는 것이 필수적입니다. DeepSORT나 ByteTrack 알고리즘을 통합하며, 이를 위해 외관 특징 임베딩 네트워크 (ReID)의 경량화 버전을 함께 배포합니다. 또한 LSTM 기반의 행동 분류를 규칙 기반 시스템과 통합하여, 더 복잡한 행동 패턴 (걷기, 뛰기, 주먹질 등)을 인식합니다. 웹 인터페이스는 멀티 뷰레이아웃 (2×2 그리드 등)으로 확장되며, 각 카메라별 독립적인 설정 (ROI, 감도 임계값 등) 조정 기능이 추가됩니다. 성능 목표는 4 대 720p 카메라 입력에서 평균 10fps 이상의 처리입니다.

1.2.3 Phase 3: NVIDIA Jetson Orin Super 이식 세 번째 단계는 검증된 시스템을 **NVIDIA Jetson Orin Super**로 이식하여 엣지 컴퓨팅 환경에서의 운영을 입증합니다. 이 단계의 핵심은 TensorRT를 활용한 모델 최적화와 GStreamer 기반 하드웨어 가속 비디오 파이프라인의 구현입니다.

TensorRT 최적화는 필수적인 변환 과정입니다. PyTorch 모델 (YOLOv8-Pose, LSTM)을 ONNX 형식으로 내보낸 후, TensorRT의 FP16 또는 INT8 양자화를 적용하여 Jetson의 Tensor Core를 활용합니다. 검색 결과에 따르면, Jetson Orin Nano에서 YOLOv8-Pose는 TensorRT 변환 후 30fps 이상의 성능을 달성할 수 있습니다. 이 과정에서 동적 배치 크기, 최적화된 플러그인, 그리고 메모리 풀기 반의 효율적인 할당 전략이 필요합니다.

GStreamer 기반 카메라 입력 파이프라인도 핵심 최적화 수단입니다. Jetson의 JetPack SDK는 nvarguscamerasrc (CSI 카메라)와 nvv4l2decoder(하드웨어 디코딩)를 제공하며, 이를 활용한 zero-copy 메모리 경로 구현으로 CPU-GPU 간 불필요한 데이터 복사를 제거합니다. RTSP 입력의 경우 rtspsrc → rtpH264Depay → nvv4l2decoder → nvvidconv → appsink 파이프라인을 구성하여 하드웨어 가속된 전처리를 달성합니다.

이 단계에서 검증해야 할 사항은 개발 환경과 동일한 정확도를 유지하면서도, Jetson의 저제한된 전력 예산 (15W~60W) 내에서 지속 가능한 처리량을 달성하는 것입니다. 열스로틀링에 대한 모니터링과, 필요 시 클럭 주파수 조정을 통한 성능-전력 트레이드 오프 관리가 포함됩니다. Docker 컨테이너화도 완료하여, 개발 환경과 타겟 환경의 불일치로 인한 “it works on my machine” 문제를 방지합니다.

1.2.4 Phase 4: Windows 네이티브 앱 전환 최종 단계는 **Python 기반 웹 앱**을 **Windows 네이티브 애플리케이션**으로 전환하여 상용 배포를 준비합니다. 이 전환의 동기는 Python 런타임의 존성 제거, 더 나은 시스템 통합, 그리고 최종 사용자에게 친숙한 설치 경험 제공입니다.

기술적 접근 방식으로는 두 가지를 고려할 수 있습니다. 첫째, **PyInstaller** 기반의 실행 파일 패키징으로 Python 런타임과 모든 종속성을 단일 .exe로 묶는 방식입니다. 이는 개발 속도가 빠르지만, 바이러스 백신의 오탐지 문제와 큰 파일 크기가 단점입니다. 둘째, **C++/C# 기반의 네이티브 재구현**으로 ONNX Runtime이나 TensorRT의 C++ API를 직접 활용하는 방식입니다. 이는 최적의 성능과 작은 배포 크기를 달성하지만, 상당한 개발 리소스가 필요합니다.

UI 프레임워크 선택으로는 **PyQt6/PySide6** 또는 **Flutter**를 검토합니다. PyQt6는 Python 코드 베이스의 재사용성이 높고, PyInstaller를 통해 단일 실행 파일을 생성할 수 있습니다. Flutter는 현대적인 크로스 플랫폼 UI와 우수한 성능을 제공하며, 향후 macOS, Linux, 모바일로의 확장성이 용이합니다. DirectShow를 통한 카메라 입력 통합과, GPU 가속을 위한 DirectML 또는 CUDA 통합도 필수적입니다. 설치 프로그램은 WiX Toolset 또는 NSIS를 활용하여 전문적인 인스톨러를 구성하며, 자동 업데이트 메커니즘도 통합합니다. 이 단계의 완료는 “정식 서비스”로의 전환을 위한 마일스톤이 됩니다.

단계	기간	핵심목표	주요산출물
Phase 1	4 주	단일카메라파이프라인검증	MPS 실시간추론, Streamlit 대시보드, 규칙 기반감지
Phase 2	4 주	멀티카메라확장및추적통합	2-4 카메라동시입력, DeepSORT, LSTM 분류
Phase 3	4 주	엣지최적화및배포	TensorRT 변환, GStreamer 파이프라인, Docker 컨테이너
Phase 4	4 주	상용네이티브앱전환	PyQt/Flutter UI, 설치형 패키지, 자동업데이트

2. 하드웨어추상화계층 (HAL)

2.1 입력장치관리모듈

2.1.1 카메라소스유형별핸들러 HAR 테스트베이스 앱은 다양한 카메라소스를 유연하게 지원해야 하며, 각 소스 유형은 고유한 초기화파라미터, 연결관리전략, 그리고 오류복구 메커니즘을 요구합니다. 이를 위해 소스 유형별 전문 핸들러를 설계하고, 공통의 CameraSource 추상인터페이스 뒤에 추상화합니다.

USB 웹캠핸들러는 로지텍 C920과 같은 표준 UVC 장치를 지원합니다. OpenCV의 cv2.VideoCapture(index) 인터페이스를 기반으로 하되, V4L2(Linux), DirectShow(Windows), AVFoundation(macOS) 등의 플랫폼별 백엔드를 자동선택하여 최적의 성능을 확보합니다. 해상도 (1920x1080부터 640x480 까지) 와 프레임레이트 (30fps, 60fps)의 동적 설정을 지원하며, cv2.CAP_PROP_BUFFERSIZE 조정을 통해 지연시간을 최소화합니다. USB 대역폭 제한으로 인해 단일 컨트롤러에 2대 이상의 고해상도 웹캠을 연결하면 성능 저하가 발생할 수 있으므로, 이 경우 해상도를 1280x720으로 낮추거나 별도의 USB 컨트롤러를 사용하는 동작 품질 적응 기능을 구현합니다.

RTSP IP 카메라핸들러는 ONVIF 프로토콜을 지원하는 네트워크 카메라를 대상으로 합니다. RTSP URL 형식은 일반적으로 rtsp://username:password@ip:port/stream_path를 따르며, 인증 정보의 안전한 관리 (환경 변수 또는 암호화된 설정 파일)가 필요합니다. OpenCV의 cv2.VideoCapture은 RTSP를 기본 지원하지만, GStreamer 백엔드를 사용하면 더 세밀한 제어가 가능합니다. rtspsrc 요소에 latency=0 옵션을 설정하여 지연시간을 최소화하거나, protocols=tcp를 지정하여 UDP 패킷 손실에 민감한 환경에서 안정성을 높일 수 있습니다. 멀티캐스트 RTSP 스트림의 경우 네트워크 부하 분산 효과를 얻을 수 있습니다.

MacBook Facetime HD 핸들러는 개발 및 테스트 단계에서의 편의성을 위해 제공됩니다. AVFoundation 프레임워크를 통해 접근하며, iOS/macOS 특화 기능인 자동 노출 조정과 초점 조정을 비활성화하여 일관된 이미지 품질을 유지할 수 있는 옵션을 제공합니다.

GStreamer 파이프라인 핸들러는 Jetson 환경에서의 최적화된 영상 입력을 담당합니다. nvarguscamerasrc(CSI 카메라), v4l2src(USB 카메라), rtspsrc(IP 카메라) 등의 소스 요소와 nvvidconv, nvjpegenc 등의 하드웨어 가속 변환 요소를 조합하여, zero-copy 메모리 관리와 최고의 처리 성능을 달성합니다. 사용자가 직접 파이프라인 문자열을 입력할 수 있는 고급 모드와, 미리 정의된 템플릿을 선택하는 간편 모드를 모두 지원합니다.

카메라유형	백엔드라이브리	주요설정파라미터	특수처리
USB 웹캠	OpenCV VideoCapture	해상도, FPS, 포맷 (YUY2/MJPEG)	V4L2/DirectShow 백엔드선택, 버퍼크기조정
RTSP IP 카메라	OpenCV + FFmpeg/ GStreamer	URL, 인증, 프로토콜 (TCP/UDP)	ONVIF 자동검색, jitterbuffer, 재연결로직
Facetime HD	AVFoundation	해상도, 자동노출비활성화	macOS 권한처리, 테스트 용최적화설정
Jetson CSI/USB/RTSP	GStreamer	파이프라인문자열또는템플릿	NVMM 메모리타입, zero-copy, 하드웨어가속

2.1.2 해상도및프레임레이트동적제어 실시간 HAR 시스템의 성능은 입력영상의 해상도와 프레임레이트에 직접적으로 영향을 받습니다. 높은 해상도는 더 정확한 포즈 추정을 가능하게 하지만 연산량이 급증하고, 높은 프레임레이트는 더 부드러운 시계열 데이터를 제공하되 처리 파이프라인에 부하를 줍니다. 따라서 상황에 따라 이러한 파라미터를 동적으로 조정할 수 있는 메커니즘이 필요합니다.

적응형 품질제어 (Adaptive Quality Control)의 기본 원리는 처리 파이프라인의 실제 처리량을 모니터링하고, 목표 프레임레이트 (예: 15fps)를 유지하기 위해 입력 품질을 조정하는 것입니다. 구체적으로 다음과 같은 알고리즘을 구현합니다: (1) 최근 N 초간의 실제 출력 FPS를 측정, (2) 목표 FPS 대비 20% 이상 낮을 경우 입력 해상도를 0.75 배로 축소, (3) 여전히 부족할 경우 0.5 배로 추가 축소, (4) 반대로 여유가 있을 경우 점진적 품질 복원. 이 과정에서 이미지의 종횡비를 유지하고, Lanczos 또는 Bicubic 보간을 사용한 고품질 다운스케일링을 적용합니다.

프레임스킵 전략은 해상도 축소만으로 부족한 경우에 사용됩니다. 입력 스트림에서 모든 프레임을 처리하는 대신, 예를 들어 30fps 입력에서 2 프레임마다 1 프레임만 처리 (15fps 효과)하거나, 3 프레임마다 1 프레임만 처리 (10fps 효과)하는 방식입니다. 중요한 것은 스킵된 프레임에 대해서도 객체 추적을 통해 궤적의 연속성을 유지하는 것이며, 이를 위해 Kalman 필터 기반의 예측을 활용합니다.

ROI(Region of Interest) 기반 적응 고급 전략으로, 장면의 특정 영역만 고해상도로 처리하고 나머지는 저해상도로 처리하는 방식입니다. 예를 들어 움직임이 감지된 영역만 1080p로 처리하고, 정적인 배경은 480p로 처리하여 전체 대역폭을 절약합니다. 이를 위해 배경 차분법이나 가벼운 움직임 검출을 전처리 단계에도 입니다. 웹 앱의 실험 인터페이스에서는 이러한 파라미터를 실시간 슬라이더로 조정하고, 변경 즉시 FPS와 CPU/GPU 사용률에 미치는 영향을 시각화하여 사용자가 최적의 설정을 탐색할 수 있도록 지원해야 합니다.

2.1.3 버퍼 관리 및 프레임 드롭 전략 실시간 비디오 처리에서 버퍼 관리는 지연 시간과 처리 신뢰성 사이의 근본 적트레이드오프를 담립니다. 무제한 버퍼링은 데이터 손실을 방지하지만 메모리 사용량 급증과 지연 시간 증가를 초래하며, 반대로 지나친 버퍼 제한은 처리 속도가 입력 속도를 따라 가지 못할 때 프레임 손실을 야기합니다.

삼중 버퍼 (Triple Buffering) 구조를 권장합니다. 각 카메라에 대해 (1) 캡처 버퍼 - 카메라 드라이버로부터 직접 프레임을 받는 잠금 없는 (lock-free) 원형 버퍼로, 가장 최신 프레임을 항상 덮어쓰기 (overwriting) 하여 지연 시간을 최소화, (2) 전처리 버퍼 - 크기 변환, 색상 공간 변환 (BGR → RGB), 정규화 등의 작업을 수행하며, 멀티 스레딩을 통해 여러 프레임을 병렬로 준비, (3) 추론 입력 버퍼 - 배치 처리를 지원하며, 충분한 프레임이 모이거나 타임 아웃이 발생하면 추론을 트리거합니다.

프레임 드롭 전략은 여러 수준에서 구현됩니다. 첫 번째 수준은 캡처 단계에서의 “선제적 드롭”으로, 처리 파이프라인의 상태를 모니터링하여 밀리는 경우 새 프레임 캡처를 일시적으로 스kip합니다. 두 번째 수준은 전처리 단계에서의 “선택적 드롭”으로, 빠른 움직임이 감지된 프레임만 우선 처리하고 정적인 장면은 낮은 우선 순위로 처리합니다. 세 번째 수준은 추론 단계에서의 “품질 기반 드롭”으

로, confidence score 가 낮은 결과는 후처리를 생략하고 다음 배치로 진행합니다. 이러한 다층 적드롭 전략은 시스템이 과부하 상황에서도 graceful degradation 을 달성하게 합니다.

버퍼 유형	크기 설정	드롭 정책	목적
캡처 버퍼	1-3 프레임	최신 프레임 덮어쓰기	지연 시간 최소화
전처리 큐	5-10 프레임	오래된 프레임 버리기	처리량 균형
추론 입력	동적 (배치 크기)	타임 아웃 기반 강제 실행	실시간 성보장

2.2 컴퓨팅 백엔드 선택기

2.2.1 Apple Silicon (MPS 가속) M1 Pro 및 후속 Apple Silicon 칩에서는 Metal Performance Shaders(MPS)를 통해 GPU 가속을 활용할 수 있습니다. PyTorch 2.1 이상에서는 MPS 백엔드를 네이티브로 지원하며, `torch.device("mps")` 설정으로 CUDA 와 유사한 API로 접근할 수 있습니다. 그러나 MPS의 완전한 호환성은 보장되지 않습니다. 특히 YOLOv8-Pose의 경우, Ultralytics 공식 문서와 GitHub 이슈를 통해 MPS에서의 알려진 제한 사항이 문서화되어 있습니다. 특히 포즈 모델의 훈련 단계에서는 MPS 관련 버그가 보고되어 있어, 공식적으로는 CPU 사용이 권장되며, 추론 단계에서는 MPS가 정상 작동합니다.

MPS 백엔드의 실제 성능은 모델 크기와 입력 해상도에 따라 달라집니다. **YOLOv8n-pose**의 경우, M1 Pro에서 **CPU-only 대비 2-3 배의 속도 향상**이 기대되며, 640x640 입력에서 30fps 이상의 실시간 처리가 가능합니다. 그러나 MPS는 NVIDIA CUDA에 비해 메모리 관리와 커널 최적화 면에서 여전히 성숙도가 떨어지므로, 대형 모델 (yolov8l/x-pose)이나 높은 배치 크기에서는 메모리 부족 오류가 발생할 수 있습니다. 따라서 MPS 환경에서는 모델 크기를 n 또는 s로 제한하고, 배치 크기를 1-4로 유지하는 것이 안정적입니다.

MPS 사용 시 주의 할 점은 **메모리 조각화 (fragmentation)** 문제입니다. 장시간 실행 시 MPS 메모리풀이 조각화되어 가용 메모리가 줄어드는 현상이 보고되었으므로, 주기적인 `torch.mps.empty_cache()` 호출을 고려해야 합니다. 또한 MPS는 FP16(BFloat16) 연산을 지원하므로, `model.half()`를 통해 반정밀도 모드로 실행하면 메모리 사용량과 처리 속도를 추가로 최적화 할 수 있습니다.

2.2.2 NVIDIA CUDA/TensorRT NVIDIA GPU 환경에서는 CUDA와 TensorRT를 통한 최적화가 성능 극대화의 핵심입니다. CUDA는 PyTorch의 기본 GPU 백엔드로, `torch.device("cuda")` 설정으로 즉시 활용 가능합니다. YOLOv8-Pose의 경우, RTX 3060 이상의 GPU에서는 640x640 입력에서 100fps 이상의 처리가 가능하며, 다중 카메라 환경에서도 안정적인 실시간 성능을 제공합니다. CUDA의 주요 이점은 성숙한 생태계와 풍부한 최적화 도구입니다. NVIDIA Nsight Systems를 통해 커널 수준의 프로파일링이 가능하고, Automatic Mixed Precision(AMP)을 통해 FP16/FP32 혼합 연산으로 속도와 정확도의 균형을 조절할 수 있습니다.

TensorRT는 추론 전용 최적화 도구로, PyTorch 모델을 ONNX로 변환한 후 TensorRT 엔진으로 빌드하여 사용합니다. TensorRT의 최적화 기법으로는 레이어 퓨전 (layer fusion), 정밀도 교정 (calibration) 기반 INT8 양자화, 동적 텐서 메모리 관리, 커널 자동 튜닝 등이 있습니다. YOLOv8-Pose를 TensorRT로 변환하면, 동일한 하드웨어에서 2-3 배의 추가 가속이 가능하며, Jetson Orin Super에서는 필수적인 최적화 경로입니다. TensorRT 빌드 시 고려 할 파라미터로는 최대 배치 크기, 최적화 레벨 (preset), 타겟 정밀도 (FP32/FP16/INT8) 등이 있으며, 애플리케이션의 지연 시간-처리량 요구 사항에 따라 튜닝 됩니다.

CUDA/TensorRT 환경에서의 주요 과제는 **버전 호환성**입니다. PyTorch, CUDA Toolkit, cuDNN, TensorRT의 버전 조합이 제한적이므로, 공식 호환성 매트릭스를 참고하여 환경을 구성해야 합니다. 또한 Windows 환경에서는 CUDA 런타임의 배포가 추가로 복잡성을 유발하므로, 정적 링크 또는 conda 환경 패키징을 고려해야 합니다.

2.2.3 CPU 풀백모드 모든 가속 하드웨어가 사용 불가능한 상황을 대비한 CPU 풀백모드는 시스템의 견고성을 위해 필수적입니다. CPU-only 실행은 속도가 느리지만, 어떤 환경에서도 동작하는 보장을 제공합니다. YOLOv8-Pose의 CPU 성능은 모델 크기에 따라 크게 달라지며, yolov8n-pose는 현대적인 멀티코어 CPU에서 10-20fps를 제공할 수 있습니다.

CPU 최적화의 핵심은 **멀티스레딩과 벡터화**입니다. OpenMP를 활용한 병렬 처리와 Intel MKL 또는 OpenBLAS를 통한 최적화된 선형 대수 연산이 기본이며, YOLOv8의 경우 NCNN이나 OpenVINO와 같은 추론 프레임워크의 도입을 검토할 수 있습니다. 그러나 CPU 모드는 실시간 처리를 보장하지 못할 수 있으므로, 명확한 성능 한계를 문서화하고 사용자에게 알려야 합니다.

CPU 모드에서의 추가 최적화로는 **모델 양자화**가 있습니다. PyTorch의 동적 양자화 (dynamic quantization)를 활용하면, 가중치를 INT8로 변환하여 메모리 대역폭을 줄이고, Intel AVX-512 등의 벡터 명령어를 활용한 연산 가속을 달성할 수 있습니다. 다만, 양자화는 약간의 정확도 저하를 동반하므로, 민감한 애플리케이션에서는 주의가 필요합니다.

2.2.4 OpenVINO (Intel 최적화) Intel CPU 및 내장 GPU에서 최적의 성능을 달성하기 위한 대안으로 Intel OpenVINO Toolkit을 고려할 수 있습니다. OpenVINO는 PyTorch 모델을 Intermediate Representation (IR) 형식으로 변환하여, Intel 하드웨어 특화 최적화를 적용합니다. 특히 최근의 Intel CPU (Core 11 세대 이상)에 탑재된 VNNI (Vector Neural Network Instructions)와 내장 GPU의 DP4A 명령어를 활용한 INT8 추론 가속이 주요 이점입니다.

OpenVINO의 활용 시나리오는 두 가지입니다. 첫째, M1 Mac에서 개발된 애플리케이션을 Intel 기반 서버로 배포할 때, 코드 변경 없이 백엔드만 교체하여 최적화를 달성할 수 있습니다. 둘째, Jetson과 함께 Intel NCS2 (Neural Compute Stick 2)를 보조 가속기로 활용하는 하이브리드 구성에서, OpenVINO는 NCS2의 Myriad VPU를 효율적으로 활용합니다. OpenVINO 2023.1 이상에서는 Ultralytics YOLOv8 모델을 직접 지원하며, ov.convert_model() API를 통해 간편하게 변환할 수 있습니다.

백엔드	지원 하드웨어	상대 성능	양자화 지원	권장 용도
MPS	Apple Silicon	1.0x (기준)	FP16	M1/M2/M3 Mac 개발
CUDA	NVIDIA GPU	1.5-3.0x	FP16, INT8	데스크톱/서버 배포
TensorRT	NVIDIA GPU/ Jetson	3.0-10.0x	FP16, INT8	고 성능 엣지 배포
OpenVINO	Intel CPU/GPU	0.3-1.5x	FP16, INT8	Intel 기반 엣지
CPU (OpenMP)	범용	0.02-0.1x	없음	풀백, 알고리즘 검증

3. 포즈 추정 알고리즘 모듈

3.1 YOLOv8-Pose 통합

3.1.1 모델 크기별 옵션 (n/s/m/l/x) YOLOv8-Pose는 다양한 컴퓨팅 환경과 정확도 요구 사항에 맞춰 5 가지 크기변형을 제공하며, 테스트 베이스 앱에서는 이러한 옵션을 런타임에 동적으로 선택할 수 있도록 하여 다양한 하드웨어와 시나리오에서의 실험을 지원합니다.

모델	파라미터수	FLOPs	COCO mAP (pose)	M1 Pro MPS FPS	Jetson Orin NX TensorRT FPS
YOLOv8n- pose	3.3M	9.2B	50.4%	45-60	120+
YOLOv8s- pose	11.6M	28.6B	60.0%	25-35	80-100
YOLOv8m- pose	26.4M	78.9B	65.0%	12-18	50-70
YOLOv8l- pose	44.4M	165.2B	68.6%	6-10	30-45
YOLOv8x- pose	68.7M	257.8B	70.3%	4-6	20-30

YOLOv8n-pose 는 3.3M 파라미터와 9.2B FLOPs 로 매우 가볍지만, COCO pose 데이터셋에서 50.4% 의 mAP 를 달성하여 기본적인 행동인식에는 충분한 정확도를 제공합니다. M1 Pro 의 MPS 가속에서는 45-60fps 를 달성할 수 있어, 30fps 입력에 대한 실시간 처리를 여유 있게 만족합니다. **YOLOv8s-pose** 는 정확도가 5.8%p 높지만 연산량이 3 배 증가하므로, GPU 가속이 확실히 보장될 때 선택합니다. 중형 이상의 모델 (m/l/x) 은 서버급 GPU 나오프라인 분석에 적합하며, 엣지 디바이스에서는 실시간성을 해칠 수 있습니다.

모델 크기 선택의 실무적 접근법은 “점진적 업그레이드” 전략입니다. 초기 개발 단계에서는 yolov8n-pose 로 빠른 프로토 타이핑을 수행하고, 실제 배포 환경에서의 정확도 요구 사항이 명확해지면 yolov8s 또는 yolov8m 로 전환합니다. 웹 앱의 실험 인터페이스에서는 사용자가 실시간으로 모델 크기를 전환할 수 있게 하여, 특정 환경에서의 성능-정확도 특성을 직접 경험하게 합니다.

3.1.2 MPS/CUDA/CPU 디바이스 자동 선택 HAL 을 통한 디바이스 자동 선택은 사용자의 수동 개입 없이 최적의 실행 환경을 보장합니다. 선택 로직은 다음 우선 순위를 따릅니다: (1) 사용자 명시적 지정, (2) CUDA 사용 가능 시 CUDA, (3) MPS 사용 가능 시 MPS, (4) 그 외 CPU.

자동 선택 로직은 PyTorch 의 백엔드 사용성 API 를 활용합니다:

```
import torch

def select_optimal_device(preferred=None):
    if preferred:
        return torch.device(preferred)
    if torch.cuda.is_available():
        return torch.device("cuda")
    elif torch.backends.mps.is_available():
        # YOLOv8-Pose 추론에 한해 MPS 사용, 훈련은 CPU로 폴백
        return torch.device("mps")
    else:
        return torch.device("cpu")
```

그러나 MPS 의 경우, 앞서 언급된 바와 같이 YOLOv8-Pose 훈련 시 알려진 버그가 있으므로, 모드별로 다른 선택으로 직이 필요합니다. 추론 모드에서는 MPS 를 허용하되, 훈련 모드 또는 안정성이 우선되는 상황에서는 CPU 로 강제 폴백하는 것이 바람직합니다. 또한 MPS 가용성 확인 시 `torch.backends.mps.is_built()` 와 `torch.backends.mps.is_available()` 를 모두 체크하여, PyTorch 빌드에 MPS 가 포함되어 있고 런타임에 MPS-enabled 디바이스가 존재하는지 확인해야 합니다.

동적 오프로딩 도 고급 기능으로 고려할 수 있습니다. 멀티 카메라 환경에서 특정 카메라의 처리가 지연될 경우, 해당 스트림만 CPU 로 오프로딩 하여 다른 스트림의 GPU 처리를 보장하는 방식입니다. 이는 우선순위 기반 스케줄링과 결합하여, critical 한 카메라 (예: 출입구 감시) 는 GPU 로, 덜 중요한 카메라는 CPU 로 처리하는 차등 서비스를 구현할 수 있습니다.

3.1.3 ONNX 변환 및 TensorRT 최적화 프로덕션 배포를 위한 핵심 단계는 PyTorch 모델을 ONNX (Open Neural Network Exchange) 형식으로 변환하고, NVIDIA 환경에서는 TensorRT 로 추가 최적화하는 것입니다. ONNX 변환은 Ultralytics에서 기본 제공하는 `export()` 메서드를 통해 수행됩니다:

```
from ultralytics import YOLO

model = YOLO("yolov8n-pose.pt")
model.export(format="onnx", dynamic=True, simplify=True)
```

`dynamic=True` 는 배치 크기와 입력 해상도를 런타임에 조정할 수 있게 하며, `simplify=True` 는 ONNX Simplifier 를 통해 불필요한 연산을 제거합니다. 생성된 ONNX 모델은 Netron 등 의 도구로 시각화하여, 변환 과정에서의 오류를 사전에 감지할 수 있습니다.

TensorRT 최적화 는 ONNX 모델을 입력으로 받아, 타겟 GPU 아키텍처에 특화된 최적화된 엔진을 빌드합니다. TensorRT 의 Python API 를 활용한 빌드 스크립트 예시는 다음과 같습니다:

```
import tensorrt as trt

logger = trt.Logger(trt.Logger.INFO)
builder = trt.Builder(logger)
network = builder.create_network(1 << int(trt.NetworkDefinitionCreationFlag.EXPLICIT_BATCH))
parser = trt.OnnxParser(network, logger)

with open("yolov8n-pose.onnx", "rb") as f:
    parser.parse(f.read())

config = builder.create_builder_config()
config.max_workspace_size = 1 << 30 # 1GB
config.set_flag(trt.BuilderFlag.FP16) # FP16 최적화

engine = builder.build_engine(network, config)
with open("yolov8n-pose.engine", "wb") as f:
    f.write(engine.serialize())
```

TensorRT 빌드 시 고려할 주요 파라미터로는: (1) `max_workspace_size` - 최적화 알고리즘을 위한 임시 메모리, (2) precision mode - FP32/FP16/INT8 중 선택, (3) `max_batch_size` - 최대 배치 크기, (4) DLA (Deep Learning Accelerator) 사용 - Jetson 의 DLA 코어 활용 여부 등이 있습니다. INT8 양자화의 경우, 교정 데이터셋을 제공하여 동적 범위를 결정해야 하며, 이는 정확도와 속도의 균형을 조절하는데 중요합니다.

3.1.4 17 개 키포인트 추출 및 정규화 YOLOv8-Pose 는 COCO 데이터셋 형식의 17 개인체 키포인트를 출력합니다. 각 키포인트는 $(x, y, visibility)$ 또는 $(x, y, confidence)$ 튜플로 표현되며, 좌표는 입력 이미지 크기에 대한 정규화된 값 (0-1

범위) 또는 퍽셀좌표로 제공됩니다. 17 개 키포인트의 구체적인 정의는 다음과 같습니다:

인덱스	키포인트 이름	설명	상위연결
0	nose	코	1, 2
1	left_eye	왼쪽 눈	3
2	right_eye	오른쪽 눈	4
3	left_ear	왼쪽 귀	5
4	right_ear	오른쪽 귀	6
5	left_shoulder	왼쪽 어깨	7, 11
6	right_shoulder	오른쪽 어깨	8, 12
7	left_elbow	왼쪽 팔꿈치	9
8	right_elbow	오른쪽 팔꿈치	10
9	left_wrist	왼쪽 손목	-
10	right_wrist	오른쪽 손목	-
11	left_hip	왼쪽 엉덩이	13
12	right_hip	오른쪽 엉덩이	14
13	left_knee	왼쪽 무릎	15
14	right_knee	오른쪽 무릎	16
15	left_ankle	왼쪽 발목	-
16	right_ankle	오른쪽 발목	-

키포인트 추출 후의 정규화는 행동 인식 모델의 입력으로 사용하기 위해 필수적입니다. 정규화의 목적은 (1) 카메라 거리와 인물 크기의 변화에 불변한 표현, (2) 좌표계의 일관성, (3) 수치 안정성을 입니다. 주요 정규화 기법으로는:

1. 절대 좌표 정규화: 모든 키포인트를 척추 중심점 (hip center, 키포인트 11과 12의 중점)을 원점으로 이동시킵니다. 이는 카메라 위치나 피사체의 이동에 더 강건한 특징을 생성합니다.
2. 스케일 정규화: 어깨-엉덩이 거리 또는 키추정치로 모든 좌표를 나누어, 개인 별 체형 차이를 보정합니다.
3. 각도 기반 특징: 관절 간의 각도를 계산하여, 방향 불변성을 확보합니다. 예를 들어, 팔꿈치 각도는 어깨-팔꿈치-손목 세 점으로 계산됩니다.

정규화된 키포인트는 34 차원 벡터 (17 개 × 2)로 표현되며, 이는 LSTM이나 ST-GCN의 입력으로 직접 사용됩니다. 신뢰도 점수가 낮은 키포인트 (예: occlusion으로 인해 감지되지 않은 경우)는 특수 값 (예: -1)으로 마스킹하거나, 보간법으로 추정합니다.

3.2 MediaPipe Pose 대안

3.2.1 BlazePose 아키텍처 활용 Google의 MediaPipe Pose는 BlazePose라는 경량 CNN 아키텍처를 기반으로 한 대안적인 포즈 추정 솔루션입니다. YOLOv8-Pose와 비교하여 단일 인물 추정에 특화되어 있으며, 모델 크기가 훨씬 작아 모바일/웹 환경에서의 실행이 용이합니다.

BlazePose의 핵심 혁신은 **heatmap-free regression** 접근법입니다. 전통적인 방식의 히트맵 기반 키포인트 검출 대신, 직접적인 좌표 회귀를 수행하여 해상도에 독립적인 일관된 정확도를 달성합니다. 이는 특히 저 해상도 입력에서의 속도와 점을 제공합니다.

다. MediaPipe 는 **33 개의 키포인트**를 출력하며, 이는 YOLOv8 의 17 개보다 상세합니다. 추가 포인트는 주로 얼굴 세부 사항 (입, 코, 눈의 세부 랜드마크) 과 손가락 뼈대에 해당합니다. 그러나 HAR 시나리오에서는 이러한 세부 사항이 오히려 노이즈가 될 수 있으므로, 17 개 COCO 호환 서브셋으로 매핑하거나 가중치를 조정하는 후처리가 필요합니다.

3.2.2 CPU-최적화 라이트 모드 MediaPipe 의 **POSE_LANDMARKS_LITE** 모델은 순수 CPU 에서도 **30fps** 이상의 실시간 처리를 목표로 합니다. 이는 GPU 가속을 사용할 수 없는 환경이나, GPU 를 다른 작업 (예: 행동 분류 신경망)에 전용하고 자활 때 유용합니다. 라이트 모드의 정확도는 풀 모델 대비 다소 저하되지만, 대부분의 기본 HAR 시나리오 (걷기, 서기, 앉기 등)에서는 여전히 충분합니다.

벤치마크에 따르면, **M1 Pro CPU**에서 MediaPipe Pose 라이트 모델은 약 **81fps** 를 달성할 수 있어, 상당한 여유 있는 실시간 성능을 제공합니다. MediaPipe 의 CPU 최적화는 Google 의 XNNPACK 라이브러리를 활용한 NEON/AVX 벡터화와 연산 그래프 최적화를 통해 달성됩니다. 라이트 모드에서는 정확도와 속도 사이의 트레이드 오프를 조정할 수 있는 파라미터 (**model_complexity**: 0/1/2) 를 제공하며, 가장 가벼운 설정 (0) 에서는 처리 속도를 극대화하고, 복잡한 설정 (2) 에서는 정확도를 우선시합니다.

CPU 플랫폼	평균 처리 시간 (ms)	최소/최대 (ms)	FPS	특이사항
Apple M1	12.3	9.8 / 18.7	~81	NPU 가속 지원
Intel i7-10700K	14.6	11.2 / 21.5	~68	AVX-512 활용
Intel i5-8250U	23.4	19.1 / 35.6	~43	저전력 모바일 CPU
AMD Ryzen 5 3600	16.8	13.5 / 24.3	~59	Zen 2 아키텍처

3.2.3 다중 인물 한계 및 대응 전략 MediaPipe Pose 의 가장 큰 실용적 한계는 공식 솔루션 이 단일 인물에 최적화되어 있다는 점입니다. 여러 사람이 동시에 화면에 나타나는 CCTV 환경에서, MediaPipe 는 주요 인물 하나만 검출하거나, 여러 인물의 키포인트가 혼합된 잘못된 결과를 출력할 수 있습니다.

이러한 한계에 대한 대응 전략으로는: (1) **YOLOv8** 과 결합한 탑다운 접근법 - 별도의 객체 검출 기로 인물 영역을 추출하고, 각 영역별로 MediaPipe Pose 를 개별 적용, (2) **MediaPipe Holistic** 솔루션과 커스텀 후처리 - 다중 인물을 추론하는 비공식 확장 활용, (3) **완전한 대체** - 다중 인물 포즈 추정에 특화된 YOLOv8-Pose 로의 전환이 있습니다. 테스트 베이스 앱에서는 이러한 대안들을 모듈로 제공하여, 사용자가 시나리오에 따라 적절한 방법을 선택할 수 있도록 합니다.

3.3 알고리즘 성능 비교 도구

3.3.1 FPS/지연 시간 벤치마크 다양한 포즈 추정 알고리즘의 실제 성능을 객관적으로 비교하기 위한 자동화된 벤치마크 도구를 구축합니다. 이 도구는 동일한 테스트 비디오 세트 (다양한 해상도, 인물 수, 동작 복잡도)에 대해 각 알고리즘을 실행하고, 초당 프레임 수 (FPS), end-to-end 지연 시간 (캡처부터 렌더링 까지), 그리고 GPU/CPU 사용률을 측정합니다.

측정은 웜업 (warm-up) 단계를 포함하여 캐시 효과를 안정화 시킨 후 수행되며, 통계적 신뢰성을 위해 여러 번 반복하여 평균과 표준 편차를 계산합니다. 벤치마크 결과는 웹 앱 인터페이스에서 실시간으로 시각화되며, 사용자가 특정 하드웨어 구성에서 어떤 알고리즘 이 최적인지 데이터 기반으로 판단할 수 있도록 지원합니다. 추가로, 장시간 안정 성 테스트 (예: 1 시간 연속 실행)를 통해 메모리 누수나 성능 저하 문제도 감지합니다.

3.3.2 정확도-속도 트레이드 오프 시각화 단순한 속도 비교를 넘어, 다양한 정확도 메트릭과 처리 속도 간의 관계를 시각화하여 사용자가 자신의 요구 사항에 최적의 지점을 선택할 수 있도록 합니다. 정확도 메트릭으로는 COCO keypoints benchmark 의 OKS (Object Keypoint Similarity), PCK (Percentage of Correct Keypoints), 그리고 사용자 정의 테스트 세트에서의 의 행동 인식 정확도를 활용합니다.

이들메트릭을 Y 축에, 처리속도 (FPS) 를 X 축에배치한 scatter plot 이나 Pareto frontier 시각화를통해, 각알고리즘변형이정확도-속도공간에서차지하는위치를직관적으로파악할수있습니다. 사용자는이시각화를통해 “내애플리케이션은 25fps 이상이필요하고, 이조건에서가장정확한알고리즘은무엇인가?” 와같은질문에답을얻을수있습니다.

4. 객체추적 및 ID 관리

4.1 추적알고리즘플러그인

4.1.1 DeepSORT (권장) DeepSORT 는 Simple Online and Realtime Tracking with a Deep Association Metric 의약자로, 객체검출과외관기반재식별 (Re-ID) 을결합한실시간다중객체추적의표준접근법입니다. HAR 테스트베이스에서의핵심역할은프레임간동일인물의일관된 ID 유지를보장하여, 시계열행동분석의기초를제공하는 것입니다.

DeepSORT 의핵심구성요소는다음과같습니다:

Kalman 필터기반모션예측: 각추적대상의바운딩박스를 $(x, y, a, h, \dot{x}, \dot{y}, \dot{a}, \dot{h})$ 의 8 차원상태벡터로모델링하며, 상수속도모션모델을가정합니다. 검출이누락된프레임에서도예측된위치로추적을유지하며, 최대 30 프레임 (1 초) 까지의 occlusion 을견딜수있습니다.

외관특징임베딩: 별도의 CNN(일반적으로 ResNet 기반) 이각검출영역으로부터 128 차원의외관특징벡터를추출합니다. 이는모션정보만으로구분이어려운 occlusion 상황에서의 ID 스위칭방지에결정적입니다. HAR 시나리오에서는포즈추정의 키포인트정보도추가특징으로활용할수있으며, 이를위한커스텀연관메트릭구현이확장지점으로预备됩니다.

형가리안알고리즘기반데이터연관: 예측된트랙과새로운검출간의비용행렬을구성하고, 최적의일대일매칭을찾습니다. 비용은 모션거리 (Mahalanobis distance) 와외관거리 (코사인유사도) 의가중조합으로정의되며, 임계값을초과하는쌍은새로운트랙으로간주됩니다.

4.1.2 ByteTrack (경량대안) ByteTrack 은 2022 년에제안된최신추적알고리즘으로, “Every Detection Box Matters” 라는철학을기반으로합니다. DeepSORT 에비해외관특징네트워크가없어메모리사용량이적고, 순수하게 검출기반으로동작하여구현이단순합니다. ByteTrack 은 Kalman 필터와형가리안알고리즘을사용하되, 연관과정에서신뢰도임계값을다단계로적용하는것이특징입니다. 낮은신뢰도의검출결과까지활용하여추적성능을향상시키며, Jetson 과같은리소스제한환경에서 DeepSORT 의외관네트워크가부담스러울때 ByteTrack 을대안으로제공합니다.

4.1.3 IOU 기반간단추적 (베이스라인) 가장단순한추적방법으로, 현재프레임의검출과이전프레임의예측간 IOU 만을 사용하여연관을결정합니다. Kalman 필터나외관특징없이순수한위치중첩에의존하므로, 빠른움직임이나 occlusion 에취약 하지만구현이매우단순하고오버헤드가거의없습니다. 이방법은 “추적없음” 옵션과의비교베이스라인으로활용되며, 추적알고리즘자체의가치를정량적으로평가하는데사용됩니다.

추적알고리즘	외관특징	Occlusion 처리	메모리사용	권장시나리오
DeepSORT	128 차원 CNN 임베딩	30 프레임예측유지	중간	복잡한 occlusion, Re-ID 필요
ByteTrack	없음 (검출신뢰도활용)	다단계임계값	낮음	리소스제한, 빠른움직임
IOU 기반	없음	없음	최소	단순한장면, 베이스라인비교

4.2 다중객체상태관리

4.2.1 객체등장/소멸감지 추적시스템은 각 객체의 생명주기를 관리하며, 새로운 객체의 등장 (새로운 검출이 기존 궤적과 연관되지 않음) 과 소멸 (일정 프레임 동안 검출되지 않음) 을 감지합니다. 등장 이벤트는 새로운 ID 할당과 함께 로그에 기록되며, 소멸이 이벤트는 해당 객체의 궤적 데이터를 최종 저장하고 메모리에서 해제합니다. 이러한 생명주기 관리는 행동 인식의 정확한 시간 범위 정의에 중요하며, “사람이 화면에 들어와서 ~ 나갈 때까지” 의 완전한 행동 시퀀스를 보장합니다.

4.2.2 ID 스위칭 방지 전략 ID 스위칭 (ID switch) 은 두 객체가 교차하거나 가까이 지나갈 때 그들의 ID 가 바뀌는 현상으로, 다중 객체 추적의 주요 오류 유형입니다. 이를 방지하기 위한 전략으로는 (1) 외관 특징의 적극적 활용 (DeepSORT), (2) 객체 체크기와 종횡비의 일관성 검사, (3) 급격한 위치 변화에 대한 물리적 제약 적용, (4) 궤적의 속도 방향 일관성 검증 등이 있습니다. 테스트 베이스 앱에서는 이러한 전략들을 파라미터화하여, 사용자가 ID 스위칭 민감도를 조절할 수 있게 합니다. ID 스위칭 발생 시 시각 경고 (색상 변화 또는 텍스트 표시) 를 제공하여 운영자의 주목을 유도합니다.

4.2.3 궤적 버퍼 및 이력 관리 각 객체의 과거 위치 이력은 행동 인식의 입력으로 활용되므로, 효율적인 궤적 버퍼 관리가 필요합니다. 고정 크기의 Ring Buffer (예: 최근 90 프레임, 3 초 분량) 를 사용하여 메모리 사용량을 제한하면서도 충분한 시간적 문맥을 유지합니다. 궤적 데이터는 (timestamp, bbox, keypoints) 투플의 시퀀스로 저장되며, 객체 소멸 시에는 선택적으로 디스크에 영구 저장하거나 비디오 클립으로 추출합니다. 행동 인식 모듈은 이 버퍼로부터 슬라이딩 윈도우를 추출하여 시퀀스 분류를 수행합니다.

5. 행동 인식 (HAR) 엔진

5.1 시퀀스 처리 파이프라인

5.1.1 슬라이딩 윈도우 버퍼 (30 프레임/1 초) 행동 인식은 단일 프레임이 아닌 시간적 문맥을 필요로 하므로, 연속된 프레임 시퀀스를 효율적으로 관리하는 버퍼 메커니즘이 핵심입니다. 슬라이딩 윈도우 방식은 고정된 크기 (예: 30 프레임, 약 1 초 @ 30fps) 의 버퍼를 유지하면서, 새 프레임이 들어올 때마다 가장 오래된 프레임을 제거하고 새 프레임을 추가합니다. 이 방식은 연속적인 행동 인식 출력력을 생성하며, 윈도우 간 중첩 (예: 50% 오버랩) 을 설정하면 시간적 해상도를 높일 수 있습니다.

버퍼의 데이터 구조는 (`batch`, `sequence`, `features`) 형태의 텐서로, 여기서 `features` 는 17 개 키포인트의 (x, y) 좌표로 구성된 34 차원 벡터입니다. 즉, 기본 설정에서 버퍼 텐서의 `shape` 는 (1, 30, 34) 가 됩니다. 슬라이딩 윈도우의 크기는 행동의 시간적 특성에 따라 조정 가능해야 하며, 빠른 행동 (예: 주먹질) 에는 작은 윈도우 (15 프레임), 느린 행동 (예: 쓰러짐) 에는 큰 윈도우 (60 프레임) 가 적합할 수 있습니다.

5.1.2 키포인트 정규화 (상대 좌표/ 중심 점 기준) 원 시키포인트 좌표는 카메라 거리, 인물 크기, 위치에 따라 절대 값이 크게 변하므로, 행동 인식 모델의 입력으로 적합하지 않습니다. 따라서 다양한 정규화 방법을 적용하여 위치와 스케일 불변성을 확보합니다:

척추 중심 점 기준 상대 좌표: 두 엉덩이 (hip) 키포인트의 중심 점을 원점 (0,0) 으로 설정하고, 모든 키포인트를 이 원점에 대해 상대적으로 재계산합니다. 이는 인물의 화면 내 위치에 관계 없이 일관된 신체 구조 표현을 제공합니다.

어깨 거리 정규화: 두 어깨 키포인트 간의 유클리드 거리를 1.0 으로 설정하는 스케일 정규화를 수행합니다. 이는 카메라로부터의 거리에 따른 인물 크기 변화를 보정하여, 동일한 행동이 가까이 또는 멀리 서 수행되어도 유사한 입력 패턴을 생성합니다.

Z-score 정규화: 키포인트 좌표의 분포를 표준 정규 분포로 변환하여, 딥러닝 모델의 수렴을 가속화합니다. 이는 학습 데이터의 통계 특성을 활용하며, 온라인 정규화 (실시간 평균/분산 추정) 도 가능합니다.

테스트 베이스 앱에서는 이러한 정규화 방법을 실시간으로 전환하여 비교할 수 있는 인터페이스를 제공하며, 각 방법의 행동 인식 정확도에 미치는 영향을 분석할 수 있게 합니다.

5.1.3 시계열 데이터 증강 제한된 학습 데이터 상황에서 모델의 일반화를 향상시키기 위한 시계열 증강 기법을 적용합니다. 공간적 증강으로는 키포인트 좌표에 대한 작은 랜덤 노이즈 추가, 수평 뒤집기 (좌우 대칭 행동의 경우), 그리고 관절 각도의 작은 왜곡이 포함됩니다. 시간적 증강으로는 시퀀스의 시간적 리샘플링 (빠르게/느리게 재생 시뮬레이션), 랜덤한 프레임 마스킹 (occlusion 시뮬레이션), 그리고 슬라이딩 윈도우의 랜덤 시프트가 있습니다. 이러한 증강은 훈련 시에만 적용되며, 추론 시에는 원본 데이터를 사용합니다.

5.2 분류알고리즘모듈

5.2.1 LSTM/GRU 기반시퀀스분류 LSTM(Long Short-Term Memory) 은 시계열데이터처리에 강한 모델로, 연속된 프레임의 관절좌표를 입력 받아 행동을 예측합니다. YOLOv8-Pose 가 추출한 좌표의 '시간적 변화 패턴'을 학습하여 "걷기", "뛰기", "쓰러짐", "싸움" 등의 행동 클래스를 분류합니다.

구체적인 아키텍처로는 **2-Layer LSTM with Hidden size 64, Dropout(0.5)**, 그리고 **Fully Connected Layer** 를 권장합니다. 입력은 30 프레임 \times 34 차원 (17 키포인트 \times 2 좌표) 의 정규화된 키포인트 시퀀스이며, 출력은 Softmax 를 통한 행동 클래스 확률 분포입니다. GRU(Gated Recurrent Unit) 는 LSTM 보다 파라미터 수가 적어 경량화된 대안으로 활용할 수 있습니다.

구성요소	설정	설명
LSTM 레이어수	2	충분한 표현력과 훈련 안정성의 균형
Hidden size	64	메모리 효율과 성능의 sweet spot
Dropout	0.5	과적합 방지
출력 활성화	Softmax	다중 클래스 분류
손실 함수	Cross-Entropy	표준 분류 손실
옵티마이저	Adam	적응적 학습률

5.2.2 ST-GCN (그래프 신경망) ST-GCN(Spatial Temporal Graph Convolutional Network) 은 관절의 움직임을 그래프 형태로 학습하여 "걷기", "뛰기", "주먹질" 등을 분류합니다. 인체의 자연스러운 관절 연결 구조를 그래프의 노드와 엣지로 모델링하여, 공간적 관계와 시간적 변화를 동시에 포착합니다. ST-GCN 은 특히 관절 간의 상호 작용이 중요한 행동 (예: 손흔들기, 포옹) 에서 LSTM 보다 우수한 성능을 보일 수 있으나, 구현 복잡도와 계산 비용이更高습니다.

5.2.3 규칙기반 폴백 (Rule-based) 간단한 수학적 규칙으로 행동을 정의하는 방식은 즉시 구현 가능하고 해석 가능성 이 높다는 장점이 있습니다. 주요 규칙 예시:

쓰러짐 (Fall Detection): 머리의 Y 좌표가 급격히 낮아지고 (예: 0.3 초내 30% 이상 감소), 몸의 가로/세로 비율이 넓어질 때 (예: aspect ratio > 2.0).

배회 (Loitering): 특정 ROI(관심 영역) 내에서 사람의 중심 좌표의 동거리가 일정 시간 동안 (예: 60 초) 임계값 미만일 때 (예: 총 이동 거리 < 2 미터).

침입 (Intrusion): 금지된 구역 (Polygon) 안에 발목 좌표가 들어왔을 때, 또는 특정 방향으로의 경계 횡단이 감지될 때.

규칙기반 방식은 딥러닝 모델의 블랙박스 특성을 보완하여, 안전-critical 한 애플리케이션에서의 신뢰성을 높입니다. 테스트 베이스 앱에서는 규칙기반과 딥러닝 기반 결과를 융합하는 하이브리드 방식도 제공합니다.

5.3 행동 클래스 정의 및 확장

5.3.1 기본 클래스 (걷기, 서기, 앉기, 뛰기) 일상적 활동을 나타내는 4 가지 기본 행동 클래스는 대부분의 HAR 시스템의 초기가 됩니다. 이들 클래스는 데이터 셋에서 충분한 샘플을 확보하기 쉽고, 상대적으로 높은 인식 정확도를 달성할 수 있습니다. 클래스 간의 구분 특성: 걷기는 주기적인 다리 움직임, 서기는 정적인 자세, 앉기는 엉덩이 높이의 급격한 감소, 뛰기는 수직 방향 가속도의 증가를 특징으로 합니다.

5.3.2 위험 행동 (쓰러짐, 싸움, 배회) 안전 모니터링 애플리케이션에서 핵심적인 위험 행동 클래스입니다. 쓰러짐은 노인들 봄에서, 싸움은 공공 안전에서, 배회는 치매 환자 감시에서 각각 중요한 이벤트입니다. 이들 클래스는 불균형 데이터 셋 문제 (정상 행동에 비해 드문 발생) 와 높은 오탐지 비용으로 인해, 더 정교한 모델 아키텍처와 임계값 튜닝이 필요합니다.

5.3.3 사용자정의행동추가인터페이스 테스트베이스앱의핵심가치중하나는새로운행동클래스의쉬운추가입니다. 웹인터페이스를통해사용자가 (1) 행동이름정의, (2) 샘플비디오업로드또는실시간녹화, (3) 자동라벨링 (기존모델의예측을초기라벨로활용) 또는수동라벨링, (4) 모델미세조정 (fine-tuning) 또는규칙정의를수행할수있도록합니다. 새로운행동은기존클래스와의혼동을최소화하기위해, 충분히구별되는특징을가져야하며, 이를위한특징중요도분석도구도제공합니다.

6. 테스트베이스웹앱기능

6.1 실시간모니터링대시보드

6.1.1 멀티뷰비디오스트리밍 (1→4 카메라확장) Streamlit 기반의웹인터페이스는초기단계에서빠른프로토타이핑을가능하게합니다. 단일카메라뷰에서시작하여, `st.columns()`와 `st.image()`를활용한 2x2 그리드레이아웃으로 4 카메라확장을지원합니다. 각카메라뷰는독립적인세션상태를유지하며, 전체화면모드와개별카메라확대/축소기능도제공합니다. WebRTC 를활용한저지연스트리밍이나, HLS/DASH 를활용한대규모배포도향후확장옵션으로고려합니다.

6.1.2 스켈레톤오버레이시각화 검출된 17 개키포인트와그연결을실시간으로영상위에렌더링합니다. COCO 데이터셋의 기본연결 (skeleton) 구조를따르며, 키포인트신뢰도에따라색상강도를조절하여시각적피드백을제공합니다. 사용자는오버레이의투명도, 선두께, 키포인트크기를실시간으로조정할수있으며, 특정관절만선택적으로표시하는필터링기능도제공합니다.

6.1.3 행동클래스및신뢰도실시간표시 각감지된객체에대해현재예측된행동클래스와신뢰도점수를바운딩박스옆에표시합니다. 높은신뢰도 (>0.8) 는녹색, 중간 (0.5-0.8) 은노란색, 낮은 (<0.5) 은빨간색으로시각화하여운영자의주목을유도합니다. 최근 N 초간의행동이력을작은스파크라인 (sparkline) 그래프로표시하여, 행동의변화추이를한눈에파악할수있게합니다.

6.2 알고리즘실험인터페이스

6.2.1 모델/파라미터 A/B 테스트 웹앱의핵심실험기능으로, 동일한입력에대해두가지알고리즘설정을병렬실행하고결과를실시간비교합니다. 예를들어, 왼쪽패널에는 YOLOv8n-pose + LSTM, 오른쪽패널에는 YOLOv8s-pose + 규칙기반을표시하여, 정확도와지연시간의트레이드오프를직접확인할수있습니다. A/B 테스트의결과 (FPS, 정확도, 메모리사용량) 는자동으로로그에기록되고, 나중에 CSV 형식으로내보내기 가능합니다.

6.2.2 슬라이딩윈도우크기조정 행동인식의시간적문맥길이를실시간으로조정하여, 짧은윈도우 (15 프레임, 0.5 초) 와긴 윈도우 (60 프레임, 2 초) 의성능차이를실험합니다. 슬라이더위젯으로 15-120 프레임범위에서자유롭게조정가능하며, 변경 즉시행동예측결과의변화를관찰할수있습니다. 이는특정행동클래스의최적윈도우크기를경험적으로결정하는데활용됩니다.

6.2.3 임계값실시간튜닝 행동감지의민감도를조절하는다양한임계값을실시간으로변경합니다. 키포인트신뢰도임계값, 행동분류확률임계값, 규칙기반행동의물리적파라미터 (예: 쓰러짐감지의 Y 좌표변화율) 등을슬라이더로조정하고, false positive 와 false negative 의변화를즉시확인합니다. 최적의임계값조합을저장하여프로필로관리할수있습니다.

6.3 코드분석및설명기능

6.3.1 알고리즘흐름시각화 현재실행중인파이프라인의데이터흐름을다이어그램으로표시합니다. 카메라입력 → 전처리 → 포즈추정 → 추적 → 행동분류 → 출력의각단계를노드로, 데이터의존성을엣지로표현하며, 각노드의현재처리량과지연시간 을실시간으로오버레이합니다. 이는병목지점식별과시스템이해에큰도움이됩니다.

6.3.2 핵심코드스니펫하이라이트 선택된알고리즘의핵심구현을 Syntax highlighting 된코드로표시합니다. 예를들어, LSTM 의 forward pass, DeepSORT 의연관로직, 또는규칙기반쓰러짐감지의수학적조건등을보여주며, 코드라인별 로간략한주석을달아이해를돕습니다. 코드는 GitHub 저장소의특정버전과연동되어, 항상최신상태를유지합니다.

6.3.3 버전별변경이력비교 동일알고리즘의두버전간의차이를 diff 형식으로표시합니다. 예를들어, YOLOv8-Pose v1.2.0 과 v1.3.0 사이의모델아키텍처변경, 하이퍼파라미터조정, 또는버그수정내역을명확히보여주며, 성능변화의원인을추 적하는데활용됩니다. Semantic versioning 의메이저/마이너/패치변경유형도함께표시합니다.

7. 데이터관리및로깅

7.1 행동이벤트로그

7.1.1 타임스탬프기반이벤트저장 모든감지된행동이벤트는정밀한타임스탬프 (밀리초단위) 와함께구조화된형식으로저장됩니다. 로그항목: 이벤트 ID, 타임스탬프, 카메라 ID, 객체 ID, 행동클래스, 신뢰도, 지속시간, 관련키포인트시퀀스의참조. 저장형식은효율적인쿼리를위해 Parquet 또는 SQLite 를사용하며, 대용량배치처리를위해 JSON Lines 도지원합니다.

7.1.2 비디오클립자동추출 행동이벤트발생전후의영상을자동으로클립으로저장하여, 나중의검증이나증거자료로활용합니다. 클립길이는이벤트유형별로설정가능 (예: 쓰러짐: 전 5 초, 후 10 초) 하며, 중복이벤트의경우클립을병합하거나우선순위에따라선택합니다. 저장형식은 H.264 로압축하여용량을최소화하고, 메타데이터는별도의 JSON 파일로함께저장합니다.

7.1.3 JSON/CSV 내보내기 수집된로그데이터를표준형식으로내보내기하여, 외부분석도구와의연동을지원합니다. JSON 은계층적구조와메타데이터보존에, CSV 는스프레드시트호환성과간단한통계분석에각각적합합니다. 내보내기시시간 범위, 카메라, 행동클래스등의필터를적용할수있습니다.

7.2 알고리즘성능메트릭

7.2.1 추론시간분포 각파이프라인단계의처리시간을상세히기록하고분석합니다. 전처리, 포즈추정, 추적, 행동분류, 후처리/렌더링의각단계별평균, 중앙값, 95th 백분위수, 최대값을계산하여, 이상치와병목지점을식별합니다. 히스토그램과박스 플롯으로시각화하며, 시간에따른추이도추적합니다.

7.2.2 정확도-재현율곡선 행동분류모델의성능을다양한임계값에서평가합니다. 검증데이터셋에대한예측결과를바탕으로, 각행동클래스별 Precision-Recall 곡선을그리고 AP(Average Precision) 를계산합니다. 이는임계값선택의객관적근거를제공하며, 클래스간불균형문제도드러냅니다.

7.2.3 혼동행렬누적 시간에따른예측결과의혼동행렬을누적하여표시합니다. 어떤행동쌍이자주혼동되는지 (예: 걷기와뛰기) 를파악하여, 모델개선의우선순위를정합니다. 혼동행렬은정규화된형태 (행합 =1) 와원수형태모두제공합니다.

7.3 데이터셋구축도구

7.3.1 자동라벨링인터페이스 기존모델의예측을초기라벨로활용하여, 수동라벨링의효율을높입니다. 사용자는예측결과를 검토하고필요시수정만수행하면되며, 수정내역은모델개선을위한피드백으로활용됩니다. 활성학습 (active learning) 기법 을적용하여, 모델이불확실한샘플을우선적으로제안하도록합니다.

7.3.2 키포인트시퀀스시각화편집기 수집된키포인트시퀀스를 2D/3D 로시각화하고, 오류를수동으로수정할수있는도구입니다. 특정프레임의키포인트위치를드래그로조정하거나, occlusion 으로인해누락된키포인트를보간법으로채울수있습니다. 수정이력은버전관리되며, 원본데이터는보존됩니다.

8. 경보및알림시스템

8.1 이벤트기반트리거

8.1.1 특정행동감지시시각적경보 위험행동 (쓰러짐, 싸움등) 감지시화면에뚜렷한시각적표시를제공합니다. 바운딩박스 의색상변경 (녹색 → 빨간색), 화면테두리깜빡임, 또는전체화면오버레이등의방식을사용하며, 심각도수준에따라다른시각 패턴을적용합니다.

8.1.2 임계값초과시알림발송 신뢰도임계값과지속시간임계값을모두충족할때만알림을발송하여, 오탐으로인한알림피로를 방지합니다. 예를들어, 쓰러짐신뢰도 > 0.8 이 3 초이상지속될때만알림을트리거합니다. 알림내용: 이벤트유형, 발생시간, 카메라위치, 스냅샷이미지, 관련비디오클립링크.

8.1.3 중복알림방지 (클다운) 동일유형의알림이너무빈번하게발송되지않도록클다운메커니즘을적용합니다. 동일객체에대한동일행동알림은최소 5 분간격으로제한하며, 심각도가높아지는경우 (예: 쓰러짐후움직임없음) 에는예외적으로즉시알림을발송합니다.

8.2 알림채널

8.2.1 웹소켓실시간푸시 Streamlit 의세션상태를활용한브라우저내실시간알림으로, 가장기본적인채널입니다. 별도의외부서비스의존없이작동하며, 데스크톱알림 API 와연동하여브라우저가백그라운드에있어도알림을표시할수있습니다.

8.2.2 이메일/Slack 연동 (확장) SMTP 를통한이메일발송과 Slack Webhook 을통한채널메시지를지원합니다. 환경변수또는설정파일로자격증명을관리하며, 심각도수준에따라다른채널로라우팅하는규칙을설정할수있습니다 (예: 경고 = 이메일, 심각 =Slack+ 이메일).

8.2.3 로컬사운드알림 데스크톱환경에서의즉각적인청각적피드백을제공합니다. 사용자지정가능한알림소리와볼륨, 그리고 “방해금지모드”에서의자동음소거기능을포함합니다.

9. 라이브러리화및배포

9.1 핵심모듈패키징

9.1.1 har_core: 포즈추정 + 추적 가장기본적인기능을담당하는패키지로, YOLOv8-Pose/MediaPipe 통합, DeepSORT/ByteTrack 추적, 그리고 HAL 의기본구현을포함합니다. 의존성: PyTorch, OpenCV, NumPy, Ultralytics. 최소한의기능만으로빠른설치와임포트를목표로합니다.

9.1.2 har_action: 행동분류엔진 LSTM, ST-GCN, 규칙기반분류기를통합한패키지로, 시퀀스처리파이프라인과 훈련/추론유ти리티를포함합니다. 의존성: PyTorch, scikit-learn(평가메트릭). har_core의출력을입력으로받아동작하며, 독립적으로업데이트가능합니다.

9.1.3 har_pipeline: end-to-end 파이프라인 전체워크플로우를 orchestration 하는고수준패키지로, 멀티카메라관리, 설정로딩, 로깅, 알림등의기능을포함합니다. 의존성: har_core, har_action, Streamlit(선택적). 빠른시작을위한사전설정파이프라인템플릿을제공합니다.

9.2 플랫폼별빌드

9.2.1 PyPI 패키지 (pip install) 표준 Python 패키지관리자를통한배포로, pip install har-core har-action har-pipeline 명령으로설치가능합니다. 플랫폼별 wheel 파일을미리빌드하여, 사용자의로컬컴파일없이빠른설치를지원합니다. 의존성성버전은엄격하게고정하여, 환경불일치문제를방지합니다.

9.2.2 Docker 컨테이너 (Jetson 호환) NVIDIA Container Toolkit 을활용한 GPU 가속컨테이너로, Jetson 용으로는 l4t-base 이미지를기반으로합니다. 개발환경과배포환경의완전한일치를보장하며, Kubernetes 와같은오케스트레이션플랫폼과의통합도용이합니다. 멀티스테이지빌드를활용하여최종이미지크기를최소화합니다.

9.2.3 Windows 실행파일 (PyInstaller) 단일파일실행가능한배포물로, Python 런타임과모든의존성을포함합니다. --onefile 옵션의단순함과 --onedir 옵션의빠른시작시간사이의트레이드오프를고려하여, 최종사용자경험에맞게선택합니다. 코드서명인증서를적용하여, Windows Defender 의오탐지를방지합니다.

9.3 API 문서화

9.3.1 자동생성 Sphinx 문서 코드내 docstring 에서자동으로 HTML 문서를생성합니다. Google Style 또는 NumPy Style docstring 을사용하며, 타입힌트도함께문서화합니다. Read the Docs 또는 GitHub Pages 를통한무료호스팅을활용합니다.

9.3.2 Jupyter 노트북토리얼 단계별 학습을 위한 실행 가능한 예제를 제공합니다. 기본 사용법부터 고급 커스터마이징 까지, 실제 데이터와 함께 동작하는 코드를 포함합니다. Google Colab 호환성도 검증하여, 로컬 환경 없이도 체험 가능하게 합니다.

10. 개발로드맵 및 체크리스트

10.1 1 단계: 단일 카메라 프로토타입 (4 주)

주차	주요 작업	산출물	성공 기준
1 주	개발 환경 설정, HAL 초기 구현	MPS 백엔드 동작 확인	YOLOv8n-pose 30fps 이상
2 주	Streamlit 대시보드 기본 구조	실시간 영상 표시	웹 인터페이스 응답 < 100ms
3 주	규칙 기반 행동 감지 구현	쓰러짐 / 침입 감지 로직	정확도 > 80% (검증 데이터)
4 주	통합 테스트 및 문서화	Phase 1 보고서	전체 파이프라인 안정화

10.2 2 단계: 멀티 카메라 확장 (4 주)

주차	주요 작업	산출물	성공 기준
1 주	멀티 스레드 입력 파이프라인	2 카메라 동시 처리	카메라 간 지연 < 50ms
2 주	DeepSORT 통합 및 튜닝	다중 객체 추적	ID 스위칭률 < 5%
3 주	LSTM 행동 분류 통합	4 가지 기본 행동 인식	분류 정확도 > 85%
4 주	A/B 테스트 인터페이스 완성	알고리즘 비교 도구	4 카메라 10fps 안정화

10.3 3 단계: Jetson 최적화 (4 주)

주차	주요 작업	산출물	성공 기준
1 주	ONNX 변환 및 TensorRT 빌드	최적화된 엔진 파일	FP16 대비 2 배 가속
2 주	GStreamer 파이프라인 구축	하드웨어 가속 디코딩	CPU 사용률 < 30%
3 주	Docker 컨테이너화	배포 가능한 이미지	cold start < 60 초

주차	주요작업	산출물	성공기준
4 주	엣지-클라우드하이브리드설계	원격모니터링프로토콜	네트워크단절시자율동작

10.4 4 단계: Windows 네이티브전환 (4 주)

주차	주요작업	산출물	성공기준
1 주	PyQt/Flutter UI 프로토타입	네이티브앱기본구조	주요화면 1:1 구현
2 주	DirectShow 카메라통합	Windows 네이티브입력	USB/RTSP 모두지원
3 주	설치프로그램및자동업데이트	MSI/EXE 빌드	설치시간 < 5 분
4 주	최종통합테스트및릴리스	v1.0 정식버전	72 시간연속안정동작

핵심체크리스트 (개발시작시)

- Python 3.8+ 및 PyTorch 2.1+ 설치
- Ultralytics (`pip install ultralytics`)
- OpenCV (`pip install opencv-python`)
- Streamlit (`pip install streamlit`)
- 테스트용 CCTV 샘플영상확보 (YouTube 또는 자체촬영)
- M1 Pro MPS 가속확인 (`torch.backends.mps.is_available()`)
- YOLOv8n-pose 기본추론테스트

권장첫번째실험: YOLOv8-Pose 를이용해영상에서사람의관절을실시간으로추출하는코드를실행한후, 이를기반으로규칙기반쓰러짐감지로직을추가합니다. 이단계가완료되면 Phase 2 의멀티카메라확장으로진행합니다.