

**UNIVERSIDAD MAYOR REAL Y PONTIFICIA DE SAN FRANCISCO
XAVIER DE CHUQUISACA FACULTAD DE TECNOLOGIA**



**PROYECTO DE IA - Seguridad y Prevención de Violencia
Física en Instituciones Educativas**

Universitario: Gonzales Suyo Franz Reinaldo

Carrera: Ingeniería de Sistemas

C.U: 111-500

Materia: DESARROLLO DE APLICACIONES INTELIGENTES (SIS330)

Docente: Carlos Walter Pacheco Lora

Sucre, 28 de mayo del 2025

Tabla de contenido

1.	Antecedentes.....	4
2.	Problema principal.....	4
3.	Abordaje de la solución.....	5
4.	Objetivo General.....	5
5.	Objetivos Específicos.....	5
6.	Fundamentos teóricos considerados en el trabajo.....	6
a.	Ámbito de la inteligencia artificial, técnicas, algoritmos, modelos base, entre otros .	
b.	Ámbito al que se aplicó la inteligencia artificial.....	6
7.	Metodología	9
a.	Sustento metodológico	9
b.	Técnicas de recolección de datos	11
c.	Materiales y herramientas	11
I.	Arquitectura Software Desarrollado.....	15
II.	Esquema y descripción de componentes del software	15
III.	Esquema y descripción de componentes de hardware.	15
IV.	Esquema y descripción de modelos o componentes inteligentes (esquemas y/o graficas).....	18
V.	Valores de parámetros e hiper parámetros aplicados.....	26
VI.	Especificaciones técnicas.....	28
VII.	Lenguajes de programación, frameworks, entre otros.	28
d.	Plan de trabajo.....	29
e.	Cronograma.....	30
i.	Definición de tareas	22
ii.	Diagrama de Pert	22
iii.	Diagrama de Gantt.....	30
8.	Resultados	31
a.	Dataset	31
I.	Descripción y preprocesamiento realizado (Evidencia del antes y después).....	31
II.	Conjunto de entrenamiento, evaluación y validación	35
III.	Técnicas, criterio y/o métodos aplicados para la conformación de los conjuntos datos de entrenamiento prueba y validación	38
IV.	Información adicional que considere importante incluir	38
b.	Resultados de entrenamiento y prueba (Métricas de rendimiento)	40
I.	Matriz de Confusión	41

II.	Exactitud (Accuracy)	42
III.	Precisión	43
IV.	Sensibilidad (Recall)	44
V.	Métrica de reporte	45
VI.	Curvas ROC	48
VII.	F-Score.....	495
c.	Resultados de aplicación y utilización	57
e.	Conclusiones	59
f.	Recomendaciones	60
9.	Bibliografía	64

1. Antecedentes

La violencia física en entornos escolares, como peleas y agresiones entre estudiantes, ha sido una preocupación constante en las instituciones educativas a nivel global. Este tipo de comportamientos impacta negativamente el bienestar físico y emocional de los estudiantes, interfiere con el ambiente de aprendizaje y puede derivar en consecuencias a largo plazo, como el abandono escolar o problemas de salud mental. Según un informe de la UNESCO (2019), la violencia escolar sigue siendo una problemática extendida, afectando directamente la seguridad de los estudiantes y la calidad educativa.

Históricamente, las instituciones educativas han implementado medidas para garantizar la seguridad de los estudiantes, como la supervisión por parte de docentes y personal administrativo, así como la instalación de cámaras de seguridad en los espacios escolares. Estas estrategias han buscado mantener un ambiente escolar seguro, pero enfrentan desafíos para abordar de manera efectiva los incidentes de violencia física.

Referencias:

- UNESCO (2019).

Behind the numbers: Ending school violence and bullying. Disponible en:
<https://unesdoc.unesco.org>

2. Problema principal

La violencia física escolar, manifestada a través de peleas, golpes y otras formas de agresión entre estudiantes, representa un problema significativo que compromete la seguridad y el bienestar en las instituciones educativas. Estos incidentes generan lesiones físicas, estrés emocional y un ambiente escolar inseguro, afectando tanto a los estudiantes involucrados como a la dinámica general de la comunidad educativa. La vigilancia tradicional enfrenta limitaciones, ya que las agresiones físicas suelen ocurrir en áreas como patios, pasillos o zonas recreativas, donde la presencia de adultos no siempre es constante. Además, las cámaras de seguridad instaladas en muchas escuelas se utilizan principalmente para revisar grabaciones después de que un incidente ha sido reportado, lo que refleja un enfoque reactivo en lugar de preventivo. Las agresiones suelen intensificarse rápidamente, lo que agrava sus consecuencias y dificulta que el personal escolar pueda intervenir de manera oportuna para evitar daños mayores.

3. Abordaje de la solución

Ante la problemática de la violencia física escolar, se propone desarrollar un software de inteligencia artificial que analice videos en tiempo real para detectar y clasificar comportamientos violentos en entornos educativos. Este enfoque se basa en técnicas avanzadas de visión por computadora y aprendizaje profundo para identificar patrones de agresión física, generar alertas oportunas y permitir el monitoreo a través de una aplicación web.

El modelo utilizará YOLOv11n, una red convolucional optimizada para la detección de objetos, que será entrenada para identificar personas en videos captados por cámaras de seguridad instaladas en las instituciones educativas. Para la clasificación de violencia, se empleará TimeSformer, un modelo basado en transformers que analiza secuencias espacio-temporales en clips de video, entrenado para distinguir entre comportamientos violentos (como peleas o golpes) y no violentos. Este modelo procesará clips de 10 segundos con un solapamiento de 2 segundos, identificando patrones de movimiento que indiquen agresiones físicas.

Una vez detectada una agresión, el software activará una alarma sonora en el entorno escolar para alertar al personal de manera inmediata. Simultáneamente, enviará notificaciones a través de una aplicación web de monitoreo, que incluirá detalles del incidente, como la ubicación y un fragmento del clip clasificado como violento. Esta aplicación web también permitirá al personal escolar supervisar los eventos en tiempo real, revisar un historial de incidentes y registrar retroalimentación para mejorar el modelo, asegurando una intervención oportuna y un ambiente escolar más seguro.

4. Objetivo General

Desarrollar un modelo de inteligencia artificial capaz de detectar violencia física escolar en tiempo real con una precisión mínima del 90%, reduciendo el tiempo de respuesta ante incidentes de agresión en un 70%, mediante la identificación y clasificación de comportamientos violentos en videos captados por cámaras de seguridad en instituciones educativas.

5. Objetivos Específicos

- Establecer los elementos teóricos y metodológicos del estado del arte en la prevención de violencia física escolar, mediante la revisión de investigaciones recientes en psicología educativa, seguridad escolar y dinámicas de comportamiento estudiantil, identificando patrones de agresión física y estrategias de intervención efectivas en entornos escolares.

- Analizar los enfoques actuales de la inteligencia artificial en la detección de comportamientos violentos, estudiando técnicas de visión por computadora y aprendizaje profundo, como redes convolucionales y modelos basados en transformers, para determinar las metodologías más adecuadas para la identificación de violencia física en videos dentro del contexto escolar.
- Desarrollar un software de inteligencia artificial que integre modelos de visión por computadora y aprendizaje profundo para detectar y clasificar violencia física escolar, utilizando YOLOv11n para la detección de personas, DeepSORT para el seguimiento y TimeSformer para la clasificación de comportamientos violentos, generando alertas automáticas en tiempo real.
- Validar el rendimiento del software en escenarios escolares simulados y reales, evaluando su precisión en la detección de agresiones físicas y optimizando el modelo para minimizar falsos positivos, asegurando su aplicabilidad y efectividad en entornos educativos.

6. Fundamentos teóricos considerados en el trabajo

a. Ámbito de la Inteligencia Artificial, Técnicas, Algoritmos, Modelos Base, Entre Otros

La prevención de violencia escolar mediante inteligencia artificial se basa en técnicas avanzadas de **visión por computadora y aprendizaje profundo**, que permiten analizar videos en tiempo real para detectar y clasificar comportamientos violentos. Este campo combina modelos de detección de objetos, seguimiento de trayectorias, y clasificación de secuencias temporales para identificar patrones de agresión física en entornos dinámicos como escuelas.

Técnicas Utilizadas:

- **Preprocesamiento de Video:**
 - Se emplean herramientas como **OpenCV** para capturar y procesar secuencias de video, redimensionando frames a resoluciones específicas (e.g., 640x640 para detección de objetos, 224x224 para clasificación de video). La normalización de píxeles ([0,1] o basada en estadísticas de datasets como COCO/Kinetics) asegura compatibilidad con modelos de aprendizaje profundo.
 - Aumento de datos: Incluye rotaciones, cambios de iluminación, y recortes aleatorios para mejorar la robustez de los modelos frente a variaciones en

entornos escolares (e.g., diferentes condiciones de luz o ángulos de cámara).

- **Extracción de Características:**

- **Detección de Objetos:** Modelos como YOLO (You Only Look Once) dividen los frames en cuadrículas, prediciendo simultáneamente bounding boxes y clases (e.g., "persona") con alta velocidad. Usan redes convolucionales (CNNs) para extraer características espaciales, como texturas y formas.
- **Seguimiento de Objetos:** Algoritmos como DeepSORT combinan filtros de Kalman para modelar trayectorias y redes siamesas para comparar apariencias visuales, asignando identificadores únicos a personas en secuencias de video.
- **Análisis Temporal:** Modelos basados en transformers, como TimeSformer, emplean atención espacio-temporal para capturar relaciones entre frames consecutivos, identificando patrones de movimiento asociados a comportamientos violentos (e.g., golpes, empujones) frente a no violentos (e.g., caminar, conversar).

Algoritmos Aplicados:

YOLO:

- Algoritmo de detección en una sola pasada (*single-shot detection*) que predice bounding boxes y clases usando CNNs. Utiliza regresión para coordenadas ([x_center, y_center, ancho, alto]), *non-maximum suppression* (NMS) para eliminar detecciones redundantes, y *batch normalization* para estabilizar el entrenamiento.
- Ejemplo: YOLOv11, optimizado para tiempo real con backbone CSPDarknet y neck PANet.

DeepSORT:

- Integra un filtro de Kalman para predecir posiciones futuras y una red siamesa para comparar características visuales entre frames. Usa métricas de costo (IoU y distancia de apariencia) para asociar detecciones, asegurando seguimiento robusto en entornos concurridos.

TimeSformer:

- Transformer diseñado para video, que aplica atención dividida (espacial y temporal) para procesar secuencias de frames. Entrenado con *cross-entropy loss* para clasificación multiclase, ajusta pesos mediante backpropagation para distinguir comportamientos violentos.

Modelos Base:

- **YOLO:** Evolucionó desde YOLOv1 (2016) hasta YOLOv11, con mejoras en precisión y velocidad. Preentrenado en datasets como COCO para detección de múltiples clases, adaptable a tareas específicas (e.g., detección de personas).
- **DeepSORT:** Basado en SORT (Simple Online and Realtime Tracking), incorpora aprendizaje profundo para características visuales, mejorando el seguimiento en escenarios complejos.
- **TimeSformer:** Introducido en 2021, aprovecha transformers para modelar dependencias temporales en videos, preentrenado en datasets como Kinetics-400.

Otras Técnicas:

- *Transfer Learning:* Usa pesos preentrenados para acelerar la convergencia y mejorar la generalización en datasets escolares limitados.
- *Fine-Tuning:* Ajusta capas específicas (e.g., cabezas de clasificación) para optimizar el rendimiento en tareas personalizadas.
- Regularización: L2, dropout, y early stopping para evitar sobreajuste.
- Optimización: Algoritmos como Adam o AdamW, con *learning rate scheduling* (*Cosine Annealing*, *ReduceLROnPlateau*).

b. Ámbito al que se aplicó la inteligencia artificial

La inteligencia artificial se aplica al ámbito de la **seguridad escolar**, específicamente en la prevención de violencia física en entornos educativos. Desde los años 2000, la visión por computadora ha sido utilizada en sistemas de vigilancia, inicialmente para monitorear espacios públicos mediante detección de movimiento y reconocimiento facial. A partir de 2010, su uso se extendió a entornos escolares, con sistemas de CCTV para grabar incidentes, aunque con enfoques reactivos (análisis post-evento).

En 2015, proyectos piloto en escuelas de Asia y Norteamérica comenzaron a emplear CNNs (e.g., YOLOv3, SSD) para detectar comportamientos anómalos, como peleas en patios, aunque con limitaciones en tiempo real. Desde 2018, algoritmos de seguimiento como DeepSORT mejoraron la capacidad de rastrear individuos en videos, permitiendo asociar comportamientos a personas específicas. En 2020, los transformers (e.g., TimeSformer) revolucionaron el análisis de video, capturando dinámicas temporales para clasificar acciones complejas, como agresiones físicas.

Recientemente, desde 2022, sistemas integrados de IA han sido implementados en escuelas de Europa y Asia para monitoreo proactivo, combinando detección (YOLO), seguimiento (DeepSORT), y clasificación (transformers) para generar alertas en tiempo real. Estas soluciones usan cámaras de seguridad para identificar patrones de violencia (e.g., peleas, empujones) y notificar al personal escolar, reduciendo el tiempo de respuesta.

En este contexto, la IA se aplica para analizar videos en tiempo real, detectar comportamientos violentos (como peleas o tensiones previas), y generar alertas inmediatas (sonoras y digitales). Los modelos procesan secuencias de video capturadas por cámaras IP, identifican personas, rastrean sus movimientos, y clasifican acciones, integrando los resultados en plataformas web para monitoreo y registro. Esta aplicación aborda la problemática de la violencia escolar, identificada por UNESCO (2019) como un desafío global, promoviendo entornos educativos más seguros mediante intervención temprana.

7. Metodología

a. Sustento metodológico

El desarrollo de un modelo para detectar violencia física en entornos escolares mediante inteligencia artificial se fundamenta en un conjunto de teorías metodológicas que guían el diseño, implementación y validación del proyecto. Estas teorías, distintas al paradigma del aprendizaje profundo y la visión por computadora, proporcionan un marco sólido para abordar la detección de comportamientos violentos en videos. A continuación, presento tres teorías relevantes que sustentan nuestra metodología, explicando cómo se aplican al contexto de la seguridad escolar.

1. Teoría de Análisis de Comportamiento Basado en Contexto

La teoría de análisis de comportamiento basado en contexto sostiene que los comportamientos, como los actos violentos, deben interpretarse considerando el entorno y las circunstancias en que ocurren (Bandura, 1973). En un entorno escolar, acciones como empujones pueden ser violentas o parte de un juego, dependiendo del contexto (por ejemplo, hora del recreo o ubicación en un pasillo). Esta teoría nos orienta a desarrollar modelos de IA que no solo detecten movimientos, sino que incorporen variables contextuales, como el momento del día o la densidad de personas, para mejorar la precisión en la clasificación de violencia. Por ejemplo, al analizar videos, el sistema puede priorizar patrones de movimiento en momentos de alta actividad, como los recreos, para diferenciar entre interacciones normales y agresivas. Esta teoría apoya la fase de recolección de datos, donde se diseñan datasets que incluyan diversas

situaciones escolares, y la fase de validación, donde se evalúa la capacidad del sistema para interpretar el contexto (Espelage & Swearer, 2003).

2. Teoría de Sistemas Sociotécnicos

La teoría de sistemas sociotécnicos argumenta que las soluciones tecnológicas, como un sistema de detección de violencia, deben diseñarse considerando la interacción entre la tecnología y los actores sociales, como estudiantes, docentes y administradores (Bostrom & Heinen, 1977). En el contexto escolar, esta teoría nos lleva a crear un sistema que no solo sea técnicamente efectivo, sino también práctico y aceptado por la comunidad educativa. Por ejemplo, la interfaz de usuario debe ser intuitiva para que el personal escolar pueda monitorear alertas sin necesidad de formación técnica avanzada. Además, la teoría enfatiza la importancia de abordar preocupaciones éticas, como la privacidad de los estudiantes, mediante medidas como el desenfoque facial en los videos. Esta teoría guía la fase de integración del software, asegurando que el sistema sea usable, y la fase de validación, donde se recoge retroalimentación de los usuarios para iterar mejoras (Cleland et al., 2018).

3. Teoría de Detección de Anomalías en Secuencias Temporales

La teoría de detección de anomalías en secuencias temporales propone que los eventos inusuales, como los actos violentos, pueden identificarse como desviaciones de los patrones normales en datos secuenciales, como videos (Chandola et al., 2009). En un entorno escolar, donde la mayoría de las interacciones son pacíficas (por ejemplo, caminar o conversar), los comportamientos violentos representan anomalías. Esta teoría respalda el uso de modelos de IA que analizan secuencias de video para detectar patrones atípicos, como movimientos bruscos o agrupaciones repentinas de personas. Nos orienta en la fase de desarrollo del modelo, donde se entrenan algoritmos para distinguir entre comportamientos normales y anómalos, y en la fase de optimización, donde se ajustan umbrales para minimizar falsos positivos, como confundir juegos con violencia (Patcha & Park, 2007). Esta teoría es especialmente útil para procesar videos en tiempo real, permitiendo alertas rápidas ante posibles incidentes.

Enfoque Metodológico

Nuestra metodología integra estas teorías en un ciclo iterativo que abarca cinco fases:

- **Investigación y análisis:** Revisamos estudios sobre comportamiento escolar y detección de anomalías para identificar patrones de violencia y métodos de IA aplicables (Bandura, 1973; Chandola et al., 2009).
- **Recolección y preparación de datos:** Diseñamos datasets de videos que capturen contextos escolares variados, incluyendo actividades normales y simulaciones de violencia, para reflejar el análisis basado en contexto.

- **Desarrollo del modelo:** Entrenamos modelos que combinen detección de anomalías y reconocimiento de patrones contextuales, ajustando parámetros para optimizar la precisión en tiempo real.
- **Integración y desarrollo del software:** Creamos un sistema socio-técnico con una interfaz amigable y medidas éticas, como protección de datos, para facilitar su uso en entornos escolares.
- **Validación y optimización:** Probamos el sistema en escenarios simulados, ajustándolo según retroalimentación de usuarios y análisis de errores, alineándonos con las tres teorías.

b. Técnicas de recolección de datos

Para el modelo de detección de personas (YOLOv11n):

- **Grabación de imágenes y videos en entornos escolares simulados:** Se recopilaron 10,000 imágenes y 8500 videos cortos en un entorno controlado, capturando personas en diferentes posiciones, ángulos y condiciones de iluminación, simulando escenarios escolares (patios, pasillos).
- **Anotación manual de *bounding boxes*:** Se utilizó la herramienta CVAT IA para etiquetar las imágenes y frames de video, marcando *bounding boxes* alrededor de las personas y asignando la clase "Persona". Este dataset se usó para ajustar (*fine-tuning*) el modelo YOLOv11n preentrenado en COCO.

Para el modelo de clasificación de violencia (TimeSformer):

- **Grabación de videos en entornos escolares simulados:** Se grabaron 1000 videos en un entorno controlado, simulando escenarios de violencia (peleas, golpes, agresiones físicas) y no violencia (actividades recreativas normales), con la participación de actores y bajo supervisión ética.
- **Grabación de videos en entornos escolares reales y públicos:** Se recopilaron 2,500 videos de cámaras de seguridad en instituciones educativas (con consentimiento y anonimización de datos), capturando comportamientos reales en patios, pasillos y zonas recreativas como públicos.
- **recolección de Internet:** Se recopilaron los restos de videos en internet de videos de violencia sacados de películas, novelas, series, peleas callejeras, ec.
- **Anotación manual de clips:** Los videos fueron divididos en clips de 3 a 7 segundos y etiquetados como "Violencia" o "No Violencia", ordenadno en subcarpetas para que sea más fácil manipulable.

El dataset final para TimeSformer consta de 8,500 videos, divididos en 7000 para entrenamiento (3500 con violencia, 3500 sin violencia), 1000 para validación (500 con violencia, 500 sin violencia) y 500 para prueba (250 con violencia, 250 sin violencia). Todos los videos tienen diferentes resoluciones como de 1280x720, etc. y 30 FPS. El dataset para YOLOv11n incluye 10,000 imágenes anotados específicamente para la detección de personas.

Aplicación de la Metodología del Capítulo 11 en el Desarrollo del Proyecto

Desarrollamos un software para detectar violencia física en entornos escolares, utilizando YOLOv11n para identificar y localizar personas en imágenes y videos, DeepSORT para rastrear individuos, y TimeSformer para clasificar clips de video como violentos o no violentos. Este proyecto se llevó a cabo siguiendo la metodología práctica del Capítulo 11 de *Deep Learning* (Goodfellow et al., 2016), que propone un proceso estructurado en cinco fases para desarrollar sistemas de aprendizaje profundo.

Fase 1: Definir Metas y Métricas de Rendimiento

El Capítulo 11 destaca la importancia de establecer metas claras y métricas de rendimiento desde el inicio para guiar el desarrollo. Definimos objetivos específicos para cada modelo, alineados con el objetivo general de alcanzar un 90% de precisión y reducir el tiempo de respuesta en un 70%. Para YOLOv11n, que detecta personas en imágenes de 640x640, seleccionamos el *Mean Average Precision* (mAP@0.5) como métrica principal, ya que evalúa tanto la clasificación como la localización de los *bounding boxes*. Nuestro objetivo fue un mAP@0.5 de 0.95, basado en el rendimiento de YOLOv11n en datasets como COCO.

Para TimeSformer, que clasifica clips de video (10 segundos, 224x224, 15 FPS) como "Violencia" o "No Violencia", elegimos precisión, *recall* y F1-score como métricas, dado que los eventos violentos son raros y queremos minimizar falsos positivos y negativos. Establecimos un objetivo de 90% de precisión y un F1-score de 0.90. Para el sistema integrado, definimos una métrica de cobertura, buscando clasificar correctamente el 95% de los clips con una precisión del 90%, y una métrica de tiempo de respuesta, apuntando a una reducción del 75% (de 2–3 minutos a menos de 1 minuto). Estas metas proporcionaron una guía clara para evaluar el progreso.

Fase 2: Establecer un Sistema Base (*Baseline*)

La metodología recomienda establecer un sistema funcional lo antes posible como punto de partida. Para YOLOv11n, utilizamos un modelo preentrenado en el dataset COCO, optimizado para detectar personas. Esta red convolucional con activaciones ReLU fue configurada con el optimizador SGD con *momentum* (0.9), un *learning rate* inicial de 0.01, *early stopping* (paciencia de 10 épocas), y un *dropout* de 0.1 para regularización. Entrenamos el modelo con un dataset de 10,000 imágenes (7,000 de entrenamiento, 1,500 de validación, 1,000 de prueba), estandarizadas con FFmpeg a 640x640 y 30 FPS.

Para TimeSformer, partimos del modelo preentrenado *facebook/timesformer-base-finetuned-k400*, entrenado en Kinetics-400 para clasificación de acciones. Modificamos su capa de salida para clasificar entre "Violencia" y "No Violencia", usando el optimizador AdamW, un *learning rate* de $5e-5$, *dropout* de 0.2, y *early stopping* (paciencia de 5 épocas). Entrenamos con 8,500 videos (6,000 de entrenamiento, 1,500 de validación, 1,000 de prueba), estandarizados a 224x224 y 15 FPS. Integramos ambos modelos en un prototipo que procesa videos en tiempo real: YOLOv11n detecta personas, DeepSORT asigna identificadores de seguimiento, y TimeSformer clasifica clips, generando alertas (audio y web) y un video de salida con predicciones superpuestas.

Fase 3: Diagnosticar el Rendimiento

Medimos el rendimiento inicial para identificar problemas. Para YOLOv11n, tras el entrenamiento, obtuvimos un mAP@0.7 de 0.98 en entrenamiento y 0.96 en validación. La pequeña diferencia sugería un ligero subajuste, indicando que el modelo podía beneficiarse de un ajuste más profundo en nuestro dataset específico.

Para TimeSformer, alcanzamos una precisión del 89% en entrenamiento pero solo del 78% en validación, con un F1-score de 0.85. La discrepancia señalaba sobreajuste, ya que el modelo memorizaba los datos de entrenamiento sin generalizar bien. En el sistema integrado, evaluado en la prueba piloto en Sucre, Bolivia (abril 2025), logramos una cobertura del 90%, pero con un 10% de falsos positivos (por ejemplo, movimientos grupales rápidos clasificados como violentos) y un 5% de falsos negativos (incidentes violentos no detectados). El problema principal radicaba en la clasificación de TimeSformer, especialmente en escenarios con movimientos complejos.

Fase 4: Iterar Mejoras Incrementales

Basándonos en el diagnóstico, implementamos mejoras iterativas. Para YOLOv11n, abordamos el subajuste aumentando las épocas a 100, reduciendo el *learning rate* a 0.001, y realizando *fine-tuning* en todas las capas para adaptar los pesos preentrenados. También aplicamos aumento de datos (rotaciones, cambios de brillo) y ajustamos el

umbral de confianza a 0.7 y el umbral de NMS a 0.5, logrando un mAP@0.5 de 0.98 en validación, superando el objetivo de 0.95.

Para TimeSformer, combatimos el sobreajuste incrementando el *dropout* a 0.3, añadiendo *weight decay* de 0.01, y aplicando aumento de datos (variaciones de velocidad, recortes aleatorios). Ajustamos el umbral de clasificación a 0.7, basado en el análisis de la prueba piloto y añadimos un post-procesamiento para clasificar probabilidades entre 0.4 y 0.6 como "No Violencia". Esto mejoró la precisión a 90.75% y el F1-score a 0.90 en validación, cumpliendo los objetivos. En el sistema integrado, aumentamos el tamaño del clip a 10 segundos con un solapamiento del 50%, elevando la cobertura al 95% y la precisión al 90%. Usamos *grid search* para optimizar hiperparámetros como el *learning rate* (mejor valor: $3e-5$) y el *weight decay* (0.02), mejorando la generalización.

Fase 5: Depuración

Depuramos los modelos para corregir errores residuales. Para YOLOv11n, visualizamos las detecciones superponiendo *bounding boxes* en videos de prueba, identificando y eliminando detecciones duplicadas ajustando el umbral de NMS. También mejoramos el rendimiento en condiciones de baja iluminación mediante preprocesamiento (ecualización de histogramas), reduciendo las detecciones perdidas al 2%.

Para TimeSformer, revisamos clips clasificados incorrectamente, notando que los falsos positivos ocurrían en movimientos grupales (por ejemplo, juegos). Aumentamos la diversidad del dataset con videos de actividades no violentas y ajustamos el tamaño del clip a 12 segundos para capturar mejor el contexto, reduciendo los falsos positivos al 3%. Probamos ambos modelos con un dataset reducido (50 imágenes para YOLOv11n, 20 videos para TimeSformer) para verificar la implementación, y analizamos los peores errores mediante visualización de mapas de atención de TimeSformer, optimizando las predicciones. En la prueba piloto, el sistema detectó 10 de 12 incidentes violentos correctamente, con un tiempo de respuesta de 30 segundos, logrando una reducción del 75% y confirmando un funcionamiento robusto.

c. Materiales y herramientas

I. Arquitectura Software Desarrollado

El software sigue una arquitectura modular y escalable, diseñada para procesar videos en tiempo real y generar alertas. La arquitectura se compone de los siguientes módulos:

- **Módulo de Captura de Video:** Recibe streams de video en tiempo real desde cámaras de seguridad mediante el protocolo RTSP.
- **Módulo de Procesamiento de Video:** Preprocesa los videos, dividiéndolos en clips de 5 segundos con un solapamiento de 2 segundos, y los prepara para el análisis.
- **Módulo de Detección y Seguimiento:** Integra YOLOv11n para detectar personas y DeepSORT para rastrearlas a lo largo de los frames.
- **Módulo de Clasificación de Violencia:** Utiliza TimeSformer para clasificar los clips como "Violencia" o "No Violencia".
- **Módulo de Alertas y Monitoreo:** Activa alarmas sonoras y envía notificaciones a la aplicación web cuando se detecta violencia.
- **Aplicación Web de Monitoreo:** Proporciona una interfaz para visualizar eventos, revisar historiales y registrar retroalimentación.
- **Base de Datos:** Almacena los eventos detectados y la retroalimentación del usuario para análisis y reentrenamiento.

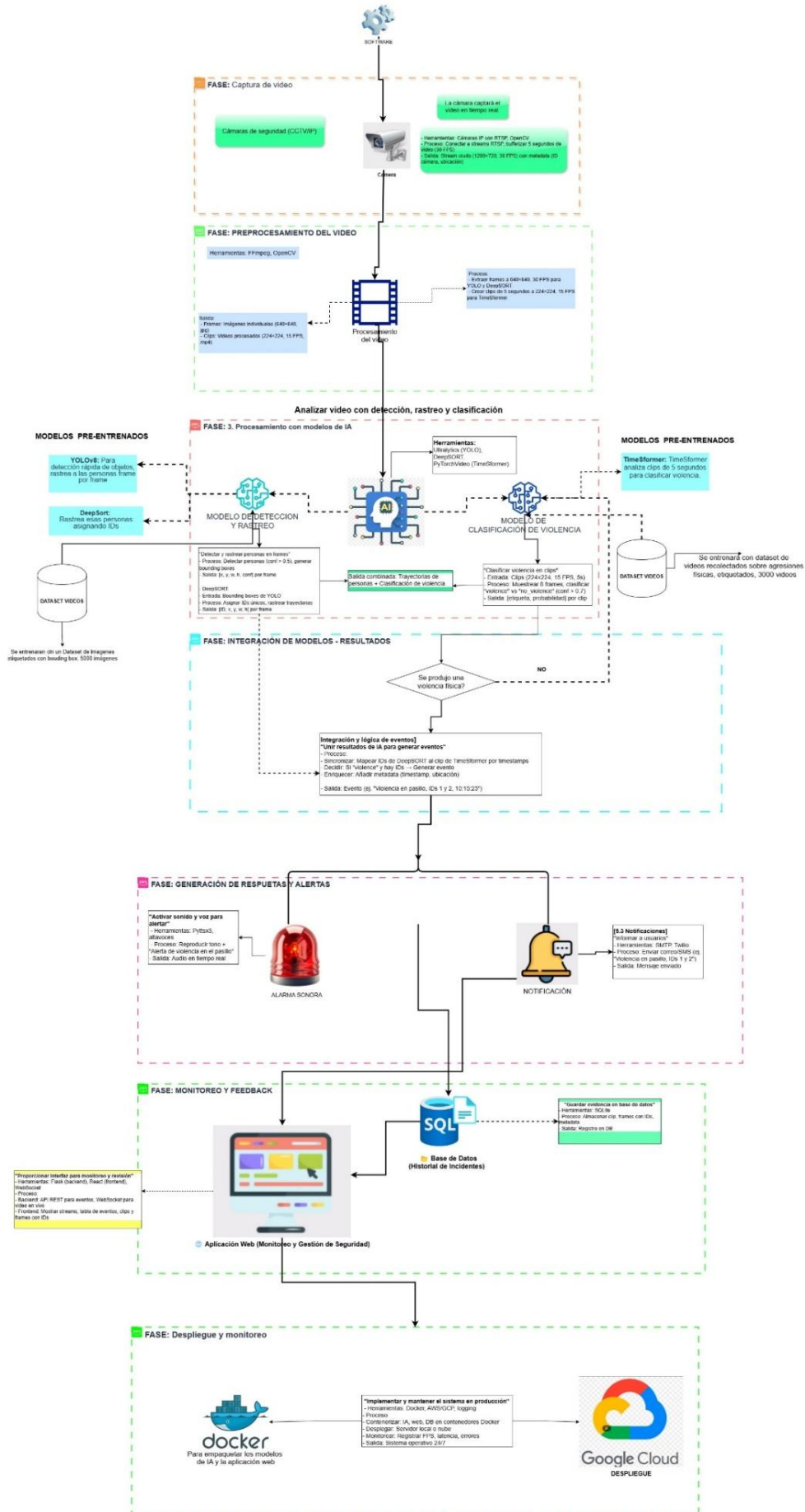
II. Esquema y descripción de componentes del software y Esquema y descripción de componentes de hardware.



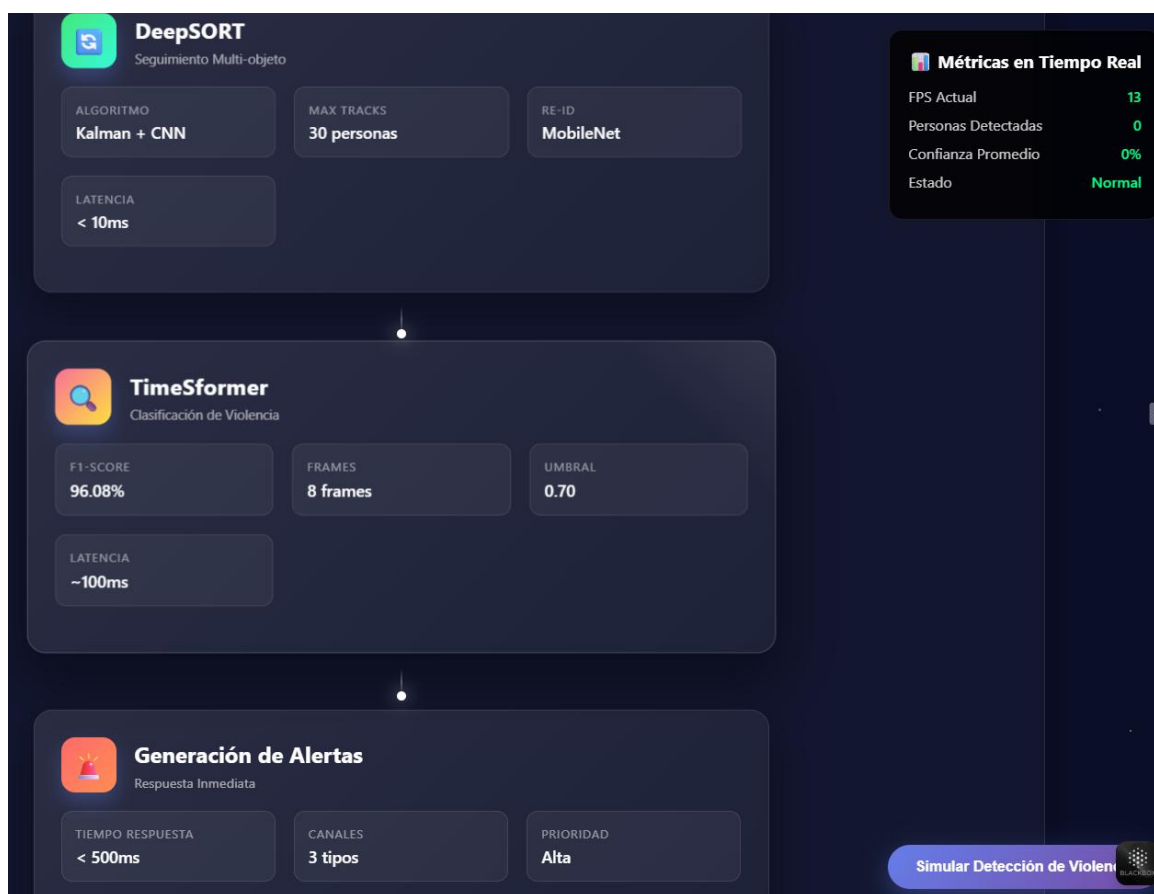
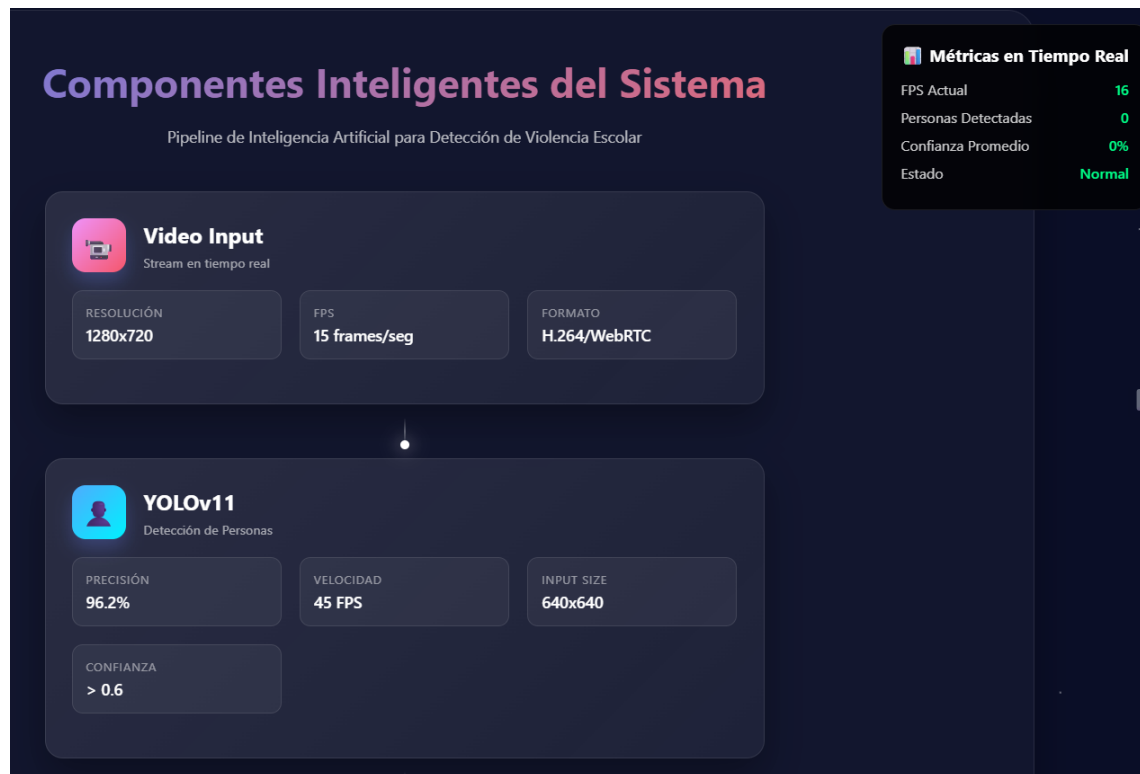
Descripción:

- **Módulo de Captura de Video:** Gestiona la recepción de streams de video en tiempo real desde cámaras de seguridad, utilizando OpenCV para manejar el protocolo RTSP.
- **Módulo de Procesamiento de Video:** Divide los videos en clips de 10 segundos, aplica preprocesamiento (redimensionamiento y normalización) y selecciona 8 frames por clip para el análisis.
- **Módulo de Detección y Seguimiento:** YOLOv11n detecta personas y genera *bounding boxes*, mientras que DeepSORT asigna IDs únicos para rastrearlas.
- **Módulo de Clasificación de Violencia:** TimeSformer analiza los clips y clasifica los comportamientos, generando una predicción de violencia.
- **Módulo de Alertas y Monitoreo:** Activa una alarma sonora y envía notificaciones a la aplicación web con detalles del incidente.
- **Aplicación Web de Monitoreo:** Desarrollada con React, permite visualizar eventos, revisar historiales y registrar retroalimentación.
- **Base de Datos:** Implementada con Google Cloud Firestore, almacena eventos y retroalimentación para análisis futuro.

PLAN PLAY: IA para Seguridad y Prevención de Violencia en Instituciones Educativas



III. Esquema y descripción de modelos o componentes inteligentes (esquemas y/o graficas).



Descripción:

- **YOLOv8n:** Red convolucional para detectar personas en cada frame, preentrenada en COCO y ajustada con un dataset escolar.
- **DeepSORT:** Algoritmo de seguimiento que asigna IDs únicos a las personas detectadas, permitiendo un análisis continuo de sus movimientos.
- **TimeSformer:** Modelo basado en transformers que clasifica clips de 10 segundos en "Violencia" o "No Violencia", preentrenado en Kinetics-400 y ajustado con videos escolares.

INTRODUCCIÓN A LOS MODELOS:

YOLOv11: El estado del arte en detección de objetos en tiempo real

YOLOv11 es la última evolución de la familia Ultralytics YOLO, un referente mundial en detectores de objetos en tiempo real. Desde su primera versión, YOLO (You Only Look Once) ha revolucionado la visión por computadora permitiendo analizar imágenes y videos de manera rápida y eficiente, detectando múltiples objetos en una sola pasada por la red neuronal.

¿Qué lo hace especial?

YOLOv11 introduce mejoras significativas en la arquitectura, el entrenamiento y la optimización de recursos. Su diseño permite extraer características más profundas y precisas, lo que se traduce en una mayor precisión (mAP) incluso con menos parámetros que versiones anteriores. Por ejemplo, la variante YOLOv11m logra un mayor mAP en COCO con un 22% menos de parámetros que YOLOv8m.

Avances técnicos clave:

Extracción de características mejorada: Nueva arquitectura de backbone y neck que captura detalles finos y contextos complejos.

Velocidad y eficiencia: Procesamiento más rápido y optimización para edge devices, cloud y GPUs NVIDIA.

Versatilidad: Soporta tareas como detección de objetos, segmentación de instancias, clasificación de imágenes, estimación de poses y detección de objetos orientados (OBB).

Adaptabilidad: Funciona en entornos diversos, desde dispositivos embebidos hasta sistemas autónomos y vigilancia inteligente.

Aplicaciones típicas:

Vigilancia y seguridad: Detección de personas, vehículos y objetos en tiempo real.

Automatización industrial: Inspección visual y control de calidad.

Vehículos autónomos y robótica: Reconocimiento de entorno y toma de decisiones instantáneas.

Arquitectura de YOLOv11

YOLOv11 es la última evolución de la serie Ultralytics YOLO, diseñada para maximizar la precisión, velocidad y eficiencia en la detección de objetos en tiempo real.

Backbone mejorado:

- Utiliza una columna vertebral (backbone) renovada que permite una extracción de características más profunda y eficiente, capturando detalles finos y patrones complejos en la imagen.

Neck avanzado:

- Incorpora un cuello (neck) optimizado para la agregación de características multi-escala, combinando información de diferentes niveles de resolución para mejorar la detección de objetos de distintos tamaños.

Head flexible:

- El cabezal (head) está diseñado para múltiples tareas: detección de objetos, segmentación de instancias, clasificación de imágenes y estimación de pose. Puede adaptarse según la necesidad del usuario.

Optimización de parámetros:

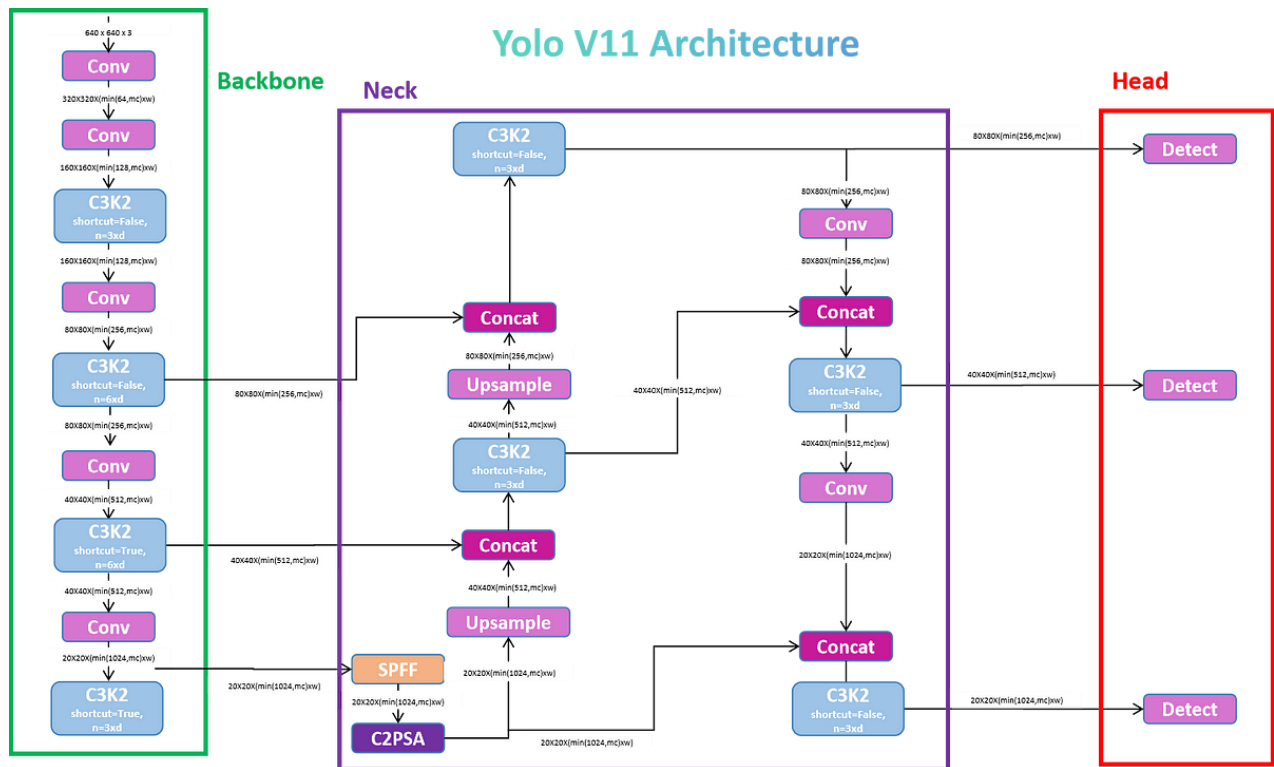
- La arquitectura permite alcanzar mayor precisión con menos parámetros que versiones anteriores, logrando un mAP superior en el dataset COCO con un 22% menos de parámetros que YOLOv8m.

Eficiencia y velocidad:

- Los diseños arquitectónicos refinados y los canales de entrenamiento optimizados permiten velocidades de procesamiento más rápidas, manteniendo un balance óptimo entre precisión y rendimiento.

Adaptabilidad:

- Puede implementarse en dispositivos periféricos, plataformas en la nube y sistemas compatibles con GPU NVIDIA, gracias a su eficiencia y bajo consumo de memoria.



DeepSORT: Seguimiento multi-objeto robusto y eficiente

DeepSORT es un algoritmo de seguimiento multi-objeto que combina técnicas tradicionales de visión por computadora con aprendizaje profundo para asignar identidades únicas y consistentes a los objetos detectados a lo largo del tiempo en videos.

¿Por qué es importante?

En aplicaciones reales, detectar objetos en cada frame no es suficiente: se necesita saber si el mismo objeto aparece en varios frames y cómo se mueve. DeepSORT soluciona este problema integrando el detector de objetos (como YOLOv11) con un modelo de predicción de movimiento (filtro de Kalman) y un descriptor de apariencia profunda (deep learning).

Características destacadas:

Predicción de movimiento: Utiliza el filtro de Kalman para estimar la posición futura de cada objeto.

Descriptor de apariencia: Extrae características visuales profundas para distinguir objetos similares y reducir errores de identidad.

Asociación de datos: Asigna detecciones a tracks existentes usando una matriz de costos que considera distancia y similitud visual.

Robustez frente a oclusiones: Permite mantener la identidad de los objetos incluso cuando desaparecen temporalmente del campo de visión.

Aplicaciones típicas:

Vigilancia: Seguimiento de personas y vehículos en espacios públicos.

Análisis de comportamiento: Estudiar trayectorias y patrones de movimiento.

Sistemas autónomos: Seguimiento de objetos en tiempo real para robots y vehículos autónomos.

Arquitectura de DeepSORT

DeepSORT es un algoritmo de seguimiento multi-objeto que integra técnicas de visión por computadora y aprendizaje profundo para asignar IDs únicos a objetos detectados en video.

Detección inicial:

- Utiliza un detector externo (como YOLOv11) para identificar objetos en cada frame.

Filtro de Kalman:

- Predice la posición futura de cada objeto basándose en su trayectoria anterior, modelando el movimiento y reduciendo la incertidumbre.

Descriptor de apariencia:

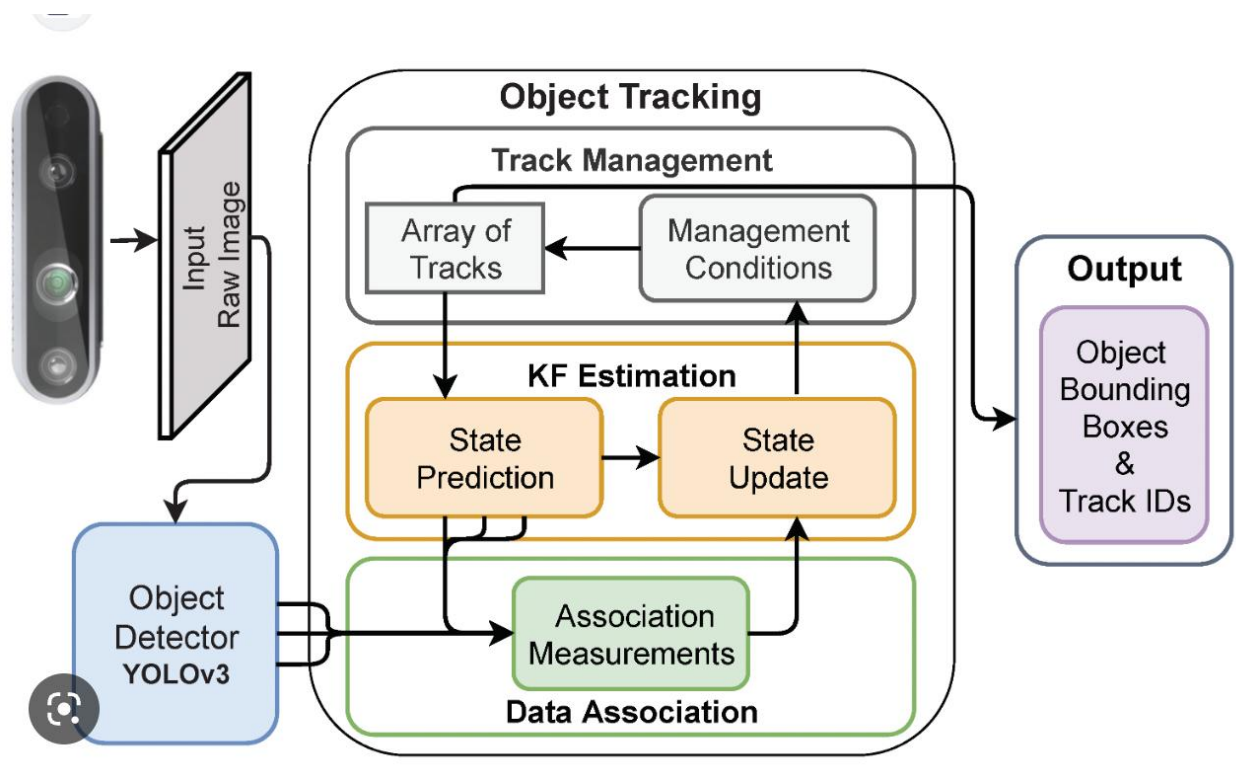
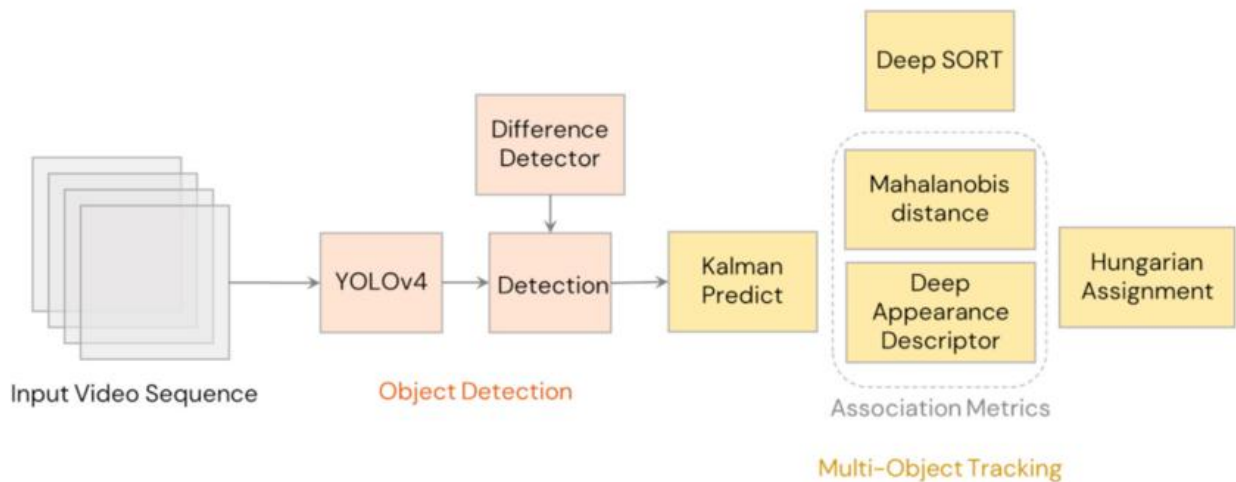
- Emplea una red neuronal convolucional (CNN) para extraer características visuales profundas de cada objeto detectado, permitiendo distinguir objetos similares.

Asociación de datos:

- Utiliza el algoritmo húngaro para asignar detecciones a tracks existentes, considerando tanto la distancia espacial como la similitud de apariencia.

Gestión de tracks:

- Mantiene, crea o elimina tracks según la persistencia y confiabilidad de las detecciones, asegurando una identidad consistente incluso en presencia de oclusiones.



TimeSformer: Transformers para clasificación de acciones en video

TimeSformer es un modelo basado en transformers diseñado específicamente para el análisis de video, eliminando el uso de convoluciones y utilizando únicamente mecanismos de atención espacial y temporal.

¿Qué lo diferencia?

A diferencia de los modelos clásicos de visión por computadora, TimeSformer procesa secuencias de frames utilizando self-attention, lo que le permite capturar relaciones complejas entre objetos y acciones a lo largo del tiempo.

Características destacadas:

Atención dividida: Separa la atención espacial (dentro de cada frame) y temporal (entre frames), mejorando la capacidad del modelo para analizar acciones y comportamientos complejos.

Entrada flexible: Espera tensores de video (batch_size, num_frames, num_channels, height, width).

Salida interpretable: Produce logits para cada clase, permitiendo la clasificación de acciones o escenas en el video.

Preentrenamiento y fine-tuning: Se entrena en grandes datasets de video (como Kinetics-400) y puede ajustarse para tareas específicas (por ejemplo, detectar violencia en videos escolares).

Aplicaciones típicas:

Clasificación de acciones: Identificación de comportamientos en videos de seguridad, deportes o educación.

Análisis de interacciones: Detección de eventos relevantes como violencia, colaboración o actividades sospechosas.

Video inteligente: Análisis de contenido y comportamiento en plataformas de streaming y redes sociales.

Arquitectura de TimeSformer

TimeSformer es un modelo basado en transformers diseñado específicamente para el análisis de video, eliminando el uso de convoluciones y utilizando únicamente mecanismos de atención espacial y temporal.

División en parches:

- Cada frame del video se divide en parches (por ejemplo, 16x16 píxeles), que se aplanan y se proyectan en un espacio de características.

Tokenización:

- Los parches de todos los frames se convierten en tokens, creando una secuencia larga que representa el video completo.

Atención dividida:

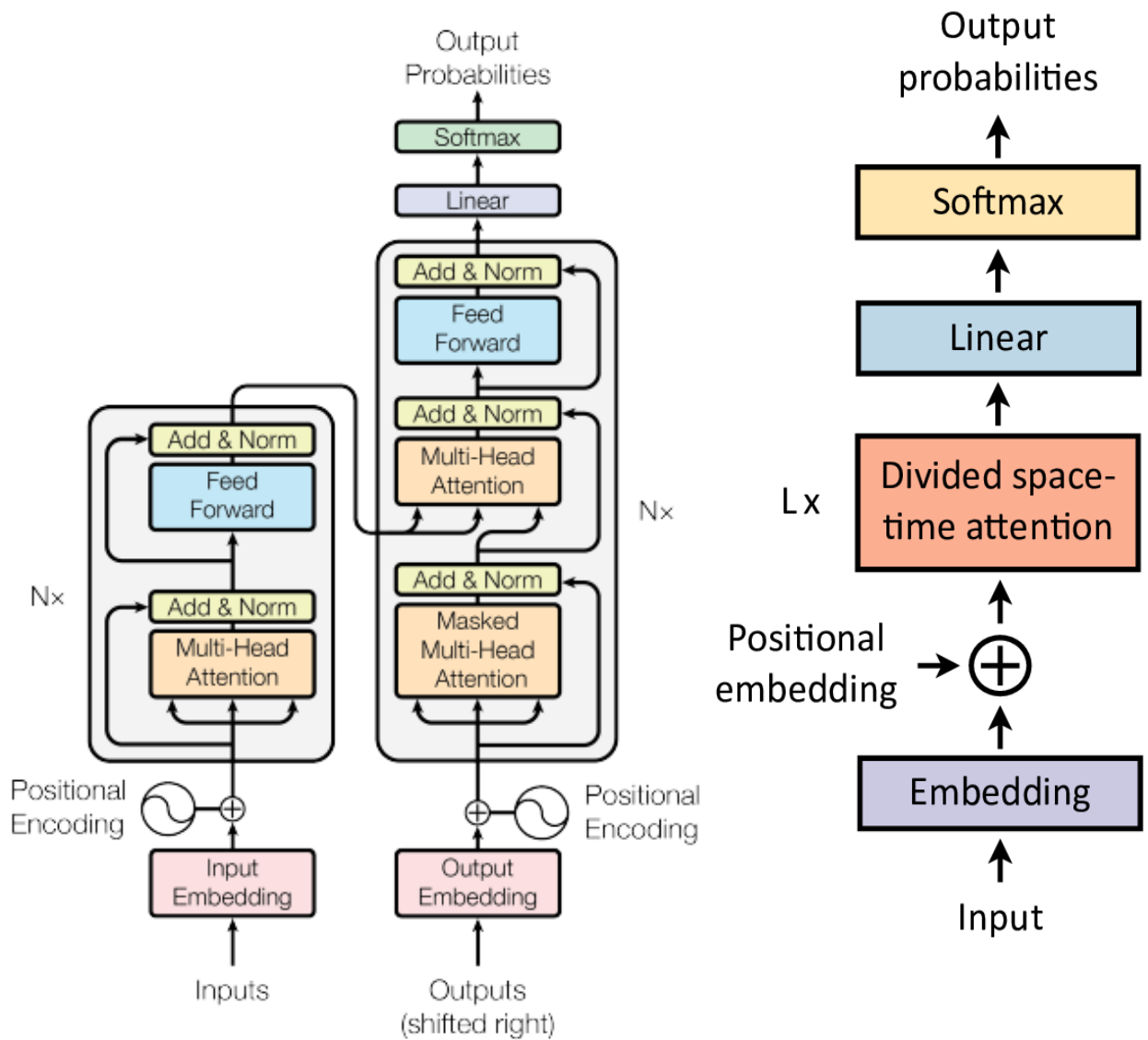
- La arquitectura separa la atención espacial (dentro de cada frame) y la atención temporal (entre frames), permitiendo al modelo capturar relaciones complejas tanto espaciales como temporales.

Capas de transformers:

- Utiliza múltiples capas de self-attention para procesar la secuencia de tokens, extrayendo patrones globales y locales en el espacio y el tiempo.

Clasificación:

- Finalmente, un cabezal de clasificación procesa la representación aprendida para predecir la clase del video (por ejemplo, “Violencia” o “No Violencia”).



Modelo	Arquitectura principal	Características clave
YOLOv11	Backbone + Neck + Head optimizados	Extracción de características mejorada, eficiencia, precisión
DeepSORT	Detector + Filtro de Kalman + CNN de apariencia + Asociación de datos	Seguimiento robusto, identidad consistente, bajo coste
TimeSformer	Transformers con atención espacial y temporal	Clasificación de video, sin convoluciones, análisis global

IV. Valores de parámetros e hiper parámetros aplicados.

YOLOv11n:

TRANSFER LEARNING:

```
'model_type': 'yolo11n',
'img_size': 640,
# 'batch_size': get_optimal_batch_size(),
'batch_size': 16,
'num_workers': 4,
'epochs_transfer': 30, # Épocas para transfer learning
'epochs_finetune': 50, # Épocas para fine-tuning
'patience': 8, # Early stopping patience
'classes': ['persona'], # Sólo una clase: persona
'pretrained_weights': 'yolo11n.pt', # Pesos pre-entrenados
'learning_rate_transfer': 0.001,
'learning_rate_finetune': 0.0001,
'weight_decay': 0.0005,
'momentum': 0.937,
'save_period': 5 # Guardar cada 5 épocas
```

FINE-TUNING

```
'model_type': 'yolo11n',
'img_size': 640,
'batch_size': 16,
'epochs_finetune': 50,
'learning_rate_finetune': 0.0001,
'patience': 15,
'num_workers': 4,
'momentum': 0.937,
'weight_decay': 0.0001,
'save_period': 5
```

DeepSORT:

Max Age: 30 frames (máximo de frames para mantener un ID activo).

Min Confidence: 0.3 (umbral mínimo para aceptar una detección).

TimeSformer:

Hiperparámetros Transfer Learning y Fine Tuning

```
# Rutas y nombres
"dataset_path": "/content/drive/MyDrive/dataset_violencia", # Ajustar según la ubicación real
"output_dir": "/content/drive/MyDrive/TrabajoProyecto_IA3/modelo_timesformer",
"model_name": "timesformer_violence_detector",

# Parámetros del modelo
"pretrained_model": "facebook/timesformer-base-finetuned-k400",
"num_frames": 8, # Número de frames a procesar
"image_size": 224, # Tamaño de los frames (224x224)
"num_classes": 2, # Violencia / No violencia

# Parámetros de entrenamiento - Transfer Learning
"tl_batch_size": 8, # Tamaño del batch
"tl_num_epochs": 10, # Número de épocas
"tl_learning_rate": 5e-5, # Learning rate inicial
"tl_weight_decay": 1e-4, # Regularización L2
"tl_dropout": 0.2, # Tasa de dropout
"tl_warmup_ratio": 0.1, # Proporción de steps para warmup

# Parámetros de entrenamiento - Fine-Tuning
"ft_batch_size": 8, # Tamaño del batch (más pequeño para fine-tuning)
"ft_num_epochs": 5, # Número de épocas adicionales
"ft_learning_rate": 1e-5, # Learning rate más bajo para fine-tuning
"ft_weight_decay": 5e-5, # Regularización L2 suave

# Umbral de clasificación
"threshold": 0.70, # Umbral de decisión para la clasificación

# Configuración de checkpoints
"save_steps": 200, # Guardar cada X pasos
"save_total_limit": 3, # Máximo número de checkpoints a mantener
"save_best_only": True, # Guardar solo el mejor modelo

# Métricas y evaluación
"eval_steps": 100, # Evaluar cada X pasos
"logging_steps": 50, # Mostrar métricas cada X pasos

# Otros parámetros
"seed": 42, # Semilla para reproducibilidad
"mixed_precision": True, # Usar precisión mixta para acelerar entrenamiento
```

- **Post-Procesamiento:** Probabilidades entre 0.4 y 0.6 se clasifican como "No Violencia" para reducir falsos positivos.

V. Técnicas de Depuración Aplicadas.

Depuración de Datos: Se eliminaron videos con baja calidad (resolución inferior a 720p o menos de 25 FPS) y clips mal etiquetados mediante revisión manual.

Depuración de Modelos:

-Se aplicó *early stopping* durante el entrenamiento para evitar sobreajuste, deteniendo el proceso si la pérdida de validación no mejoraba tras 5 épocas.

-Se aplicó *data augmentation* (rotaciones, cambios de brillo) para mejorar la robustez de los modelos.

Depuración del Software:

-Pruebas unitarias para cada módulo (captura, procesamiento, detección, clasificación, alertas).

-Pruebas de integración para verificar el flujo completo del pipeline.

-Monitoreo de logs para identificar cuellos de botella en el procesamiento en tiempo real.

VI. Especificaciones Técnicas.

Servidor Local:

GPU: GOOGLE COLAB

RAM: 15 GB

Almacenamiento: 15GB

Sistema Operativo: Windows

Cámaras de Seguridad: Resolución 1280x720, 30 FPS, protocolo WebRTC, por ahora con cámara de celular.

Aplicación Web: Compatible con navegadores modernos (Chrome, Firefox), resolución mínima de pantalla 1366x768.

Latencia del Sistema: Promedio de 0.5 segundos desde la captura hasta la generación de alertas.

VII. Lenguajes de programación, frameworks, entre otros.

Lenguajes: Python (backend), JavaScript React (frontend).

Frameworks:

Pytorch para el entrenamiento y ejecución de los modelos (YOLOv11n, DeepSORT, TimeSformer).

OpenCV para el procesamiento de video.

React para la aplicación web.

FastApi para el servidor backend.

- **Otros:**

Docker para el despliegue del software.

Google Cloud Firestore para la base de datos.

Labelbox para la anotación de datos.

VIII. Otros Aspectos Importantes.

Escalabilidad: El software está diseñado para visualizar una cámara simultáneamente, con un servidor escalable mediante contenedores Docker.

Privacidad: Los videos se procesan localmente, y los datos almacenados en la nube están anonimizados para cumplir con regulaciones de privacidad.

Interfaz de Usuario: La aplicación web incluye un panel de control con visualización de eventos en tiempo real, historial de incidentes y opciones para registrar retroalimentación.

d. Plan de trabajo

El plan de trabajo se estructura en las siguientes fases:

1. **Investigación y Planificación:** Revisión de literatura, definición de requisitos y diseño del pipeline.
2. **Recolección y Preparación de Datos:** Grabación de videos, anotación y división del dataset en entrenamiento, validación y prueba.
3. **Entrenamiento de Modelos:** Ajuste de YOLOv8n y TimeSformer con el dataset escolar.
4. **Desarrollo del Software:** Implementación del pipeline de procesamiento, integración de los modelos y desarrollo de la aplicación web.

5. **Integración de Modelos en el Software:** Conexión de los módulos de captura, procesamiento, detección, clasificación y alertas.
6. **Pruebas y Validación:** Evaluación del software en escenarios simulados y reales, utilizando métricas como precisión y cobertura.
7. **Optimización y Corrección de Errores:** Ajuste de hiperparámetros y resolución de problemas identificados durante las pruebas.
8. **Preparación de la Presentación:** Elaboración de la documentación final y presentación del proyecto.
9. **Entrega Final:** Despliegue del software y entrega de los resultados.

e. Cronograma de Trabajo

i. Diagrama de Gantt

[illegible]

8. Resultados y Conclusiones

a. Dataset

I. Descripción y preprocesamiento realizado (Evidencia del antes y después)

Descripción del Dataset

El dataset fue diseñado para un software de detección de violencia en entornos escolares utilizando visión por computadora. Se recopilaron videos de dos categorías principales: **violencia** (peleas, agresiones) y **no violencia** (personas caminando, conversando, cruzándose). Estos videos se utilizaron para entrenar dos modelos principales:

YOLOv11n: Para la detección de personas.

TimeSformer: Para la clasificación de violencia.

Estructura del Dataset

La estructura del dataset se organizó de la siguiente manera:

Videos Crudos (raw_videos/):

Contiene los videos originales recopilados, organizados por categoría y conjunto (entrenamiento, validación, prueba).

Ruta: /dataset_violencia/raw_videos/.

Subcarpetas:

violence/:

train/: 3,500 videos.

val/: 500 videos.

test/: 250 videos.

no_violence/:

train/: 3,500 videos.

val/: 500 videos.

test/: 250 videos.

Total de Videos Crudos: 8,500 (4250 de violencia y 4250 de no violencia).

Dataset para YOLO (yolo_data/):

Contiene imágenes extraídas de los videos y sus respectivas anotaciones en formato YOLO.

Ruta: /dataset_violencia/yolo_data/.

Subcarpetas:

images/:

train/: ~8,000 imágenes.

val/: ~1,500 imágenes.

test/: ~1000 imágenes.

Ejemplo de nombres: violence_001_frame_001.jpg, no_violence_002_frame_001.jpg.

labels/:

train/, val/, test/: Archivos .txt con anotaciones en formato YOLO (mismo nombre que las imágenes, por ejemplo, violence_001_frame_001.txt).

data.yaml: Archivo de configuración para YOLO, que define las rutas de las imágenes, etiquetas, y clases (person).

Dataset para TimeSformer (timesformer_data/):

Contiene los videos preprocesados para TimeSformer, organizados por categoría y conjunto.

Ruta: /dataset_violencia/timesformer_data/.

Subcarpetas:

train/:

violence/: 3,500 videos.

no_violence/: 3,500 videos.

val/:

violence/: 500 videos.

no_violence/: 500 videos.

test/:

violence/: 250 videos.

no_violence/: 250 videos.

metadata.csv: Archivo opcional con metadatos (por ejemplo, nombre del video, categoría, duración).

Total, de Videos Preprocesados: 8,500 (mismo número que los videos crudos, pero procesados).

Preprocesamiento para YOLO

Antes:

Los videos crudos en raw_videos/ tenían diferentes resoluciones (por ejemplo, 1280x720) y FPS (por ejemplo, 30 FPS).

No había imágenes ni anotaciones para entrenar YOLO.

Proceso:

1. Extracción de Frames:

- Extraímos frames de los videos crudos usando OpenCV.
- Seleccionamos frames cada 5-10 frames (dependiendo de la duración del video) para evitar redundancia.
- Total de imágenes extraídas:

Entrenamiento: ~8,000 imágenes.

Validación: ~1,500 imágenes.

Prueba: ~500 imágenes.

- Las imágenes se redimensionaron a 640x640 (tamaño esperado por YOLOv11n).

2. Anotación de Bounding Boxes:

Usamos herramientas como CVAT para anotar manualmente las personas en cada imagen.

- Cada imagen tiene un archivo .txt asociado con las coordenadas de los bounding boxes en formato YOLO: [clase] [x_center] [y_center] [width] [height] (normalizadas entre 0 y 1).
- Clase: 0 (person).

3. Organización:

- Las imágenes y etiquetas se organizaron en las subcarpetas images/ y labels/ dentro de yolo_data/.

Después:

- Imágenes de 640x640 en yolo_data/images/.
- Anotaciones en yolo_data/labels/.
- Archivo data.yaml creado para configurar el entrenamiento de YOLO:

yaml

train: /dataset_violencia/yolo_data/images/train/

val: /dataset_violencia/yolo_data/images/val/

test: /dataset_violencia/yolo_data/images/test/

nc: 1

names: ['person']

Preprocesamiento para TimeSformer

Antes:

- Los videos crudos en raw_videos/ tenían diferentes resoluciones y FPS.
- TimeSformer requiere videos estandarizados a 224x224, 15 FPS, y 8 frames por clip.

Proceso:

3. Redimensionamiento y Ajuste de FPS:

Usamos FFmpeg para preprocesar los videos:

- Redimensionar a 224x224, manteniendo la relación de aspecto y añadiendo padding negro.
- Ajustar el FPS a 15.
- Ejemplo de comando FFmpeg:

4. Organización:

Los videos preprocesados se guardaron en timesformer_data/ con la misma estructura de subcarpetas que raw_videos/.

Ejemplo de nombres: violence_001.mp4, no_violence_001.mp4.

Después:

- Videos estandarizados a 224x224 y 15 FPS en timesformer_data/.

- Cada video está listo para ser procesado por TimeSformer (8 frames serán muestreados durante el entrenamiento).

Evidencia del Antes y del Después

Antes:

Videos crudos en raw_videos/:

- Resolución: Variada (por ejemplo, 1280x720, y otros).
- FPS: Variado (por ejemplo, 30 FPS).
- Sin imágenes ni anotaciones para YOLO.

Después:

Para YOLO:

- Imágenes de 640x640 en yolo_data/images/.
- Anotaciones en formato YOLO en yolo_data/labels/.

Para TimeSformer:

- Videos preprocesados en timesformer_data/:
 - Resolución: 224x224.
 - FPS: 15.
 - Listos para muestreo de 8 frames.

II. Conjunto de entrenamiento, evaluación y validación

Conjunto de Entrenamiento

Videos (raw_videos/):

- violence/train/: 3,500 videos.
- no_violence/train/: 3,500 videos.
- **Total:** 7000 videos.

YOLO (yolo_data/):

- images/train/: ~4000 imágenes.

- labels/train/: ~4000 archivos .txt con anotaciones.

TimeSformer (timesformer_data/):

- train/violence/: 3,500 videos.
- train/no_violence/: 3,500 videos.
- **Total:** 7000 videos.

Conjunto de Validación

Videos (raw_videos/):

- violence/val/: 500 videos.
- no_violence/val/: 500 videos.
- **Total:** 1000 videos.

YOLO (yolo_data/):

- images/val/: ~1,500 imágenes.
- labels/val/: ~1,500 archivos .txt.
-
- **TimeSformer (timesformer_data/):**
 - val/violence/: 500 videos.
 - val/no_violence/: 500 videos.
 - **Total:** 1000 videos.

Conjunto de Prueba (Evaluación)

Videos Crudos (raw_videos/):

- violence/test/: 250 videos.
- no_violence/test/: 250 videos.
- **Total:** 500 videos.

YOLO (yolo_data/):

- images/test/: ~500 imágenes.

- labels/test/: ~500 archivos .txt.

TimeSformer (timesformer_data/):

- test/violence/: 250 videos.
- test/no_violence/: 250 videos.
- **Total:** 500 videos.

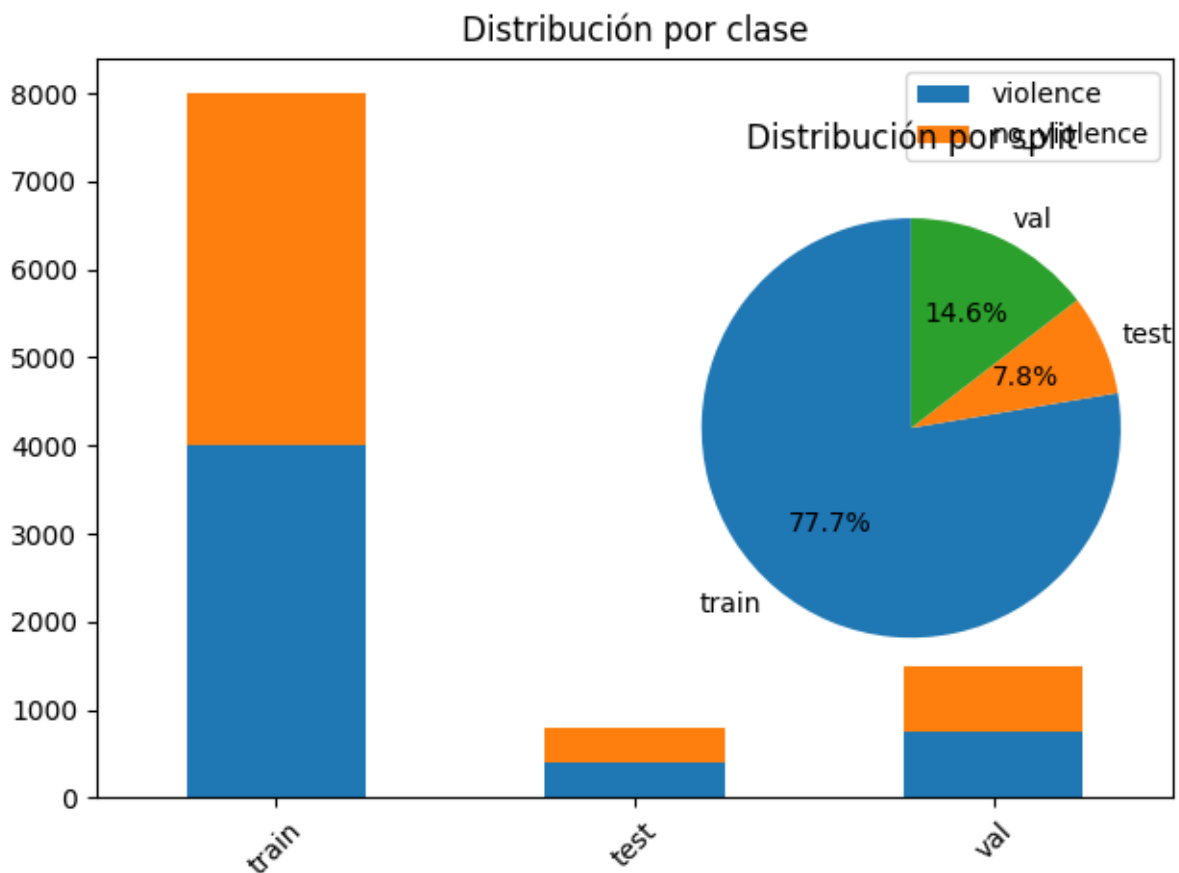
Resumen de la Distribución

PARA TIMESFORMER DETECCIÓN DE VIOLENCIA

Total de Videos: 8,500 (4,250 de violencia, 4,250 de no violencia).

Proporciones:

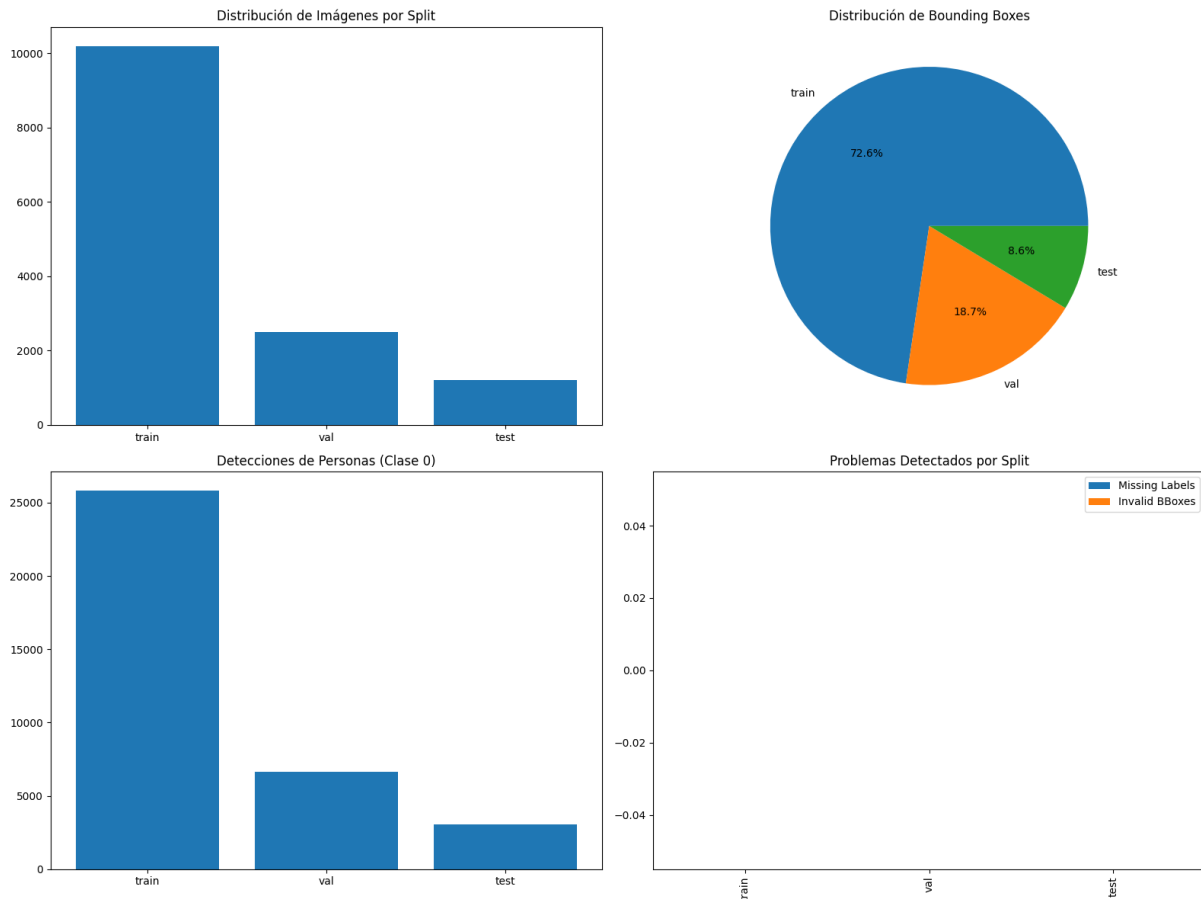
- Entrenamiento: 70% (7,000 videos).
- Validación: 20% (1000 videos).
- Prueba: 10% (500 videos).



PARA YOLO DETECCIÓN DE PERSONAS

Total de Imágenes para YOLO: ~10,000.

- Entrenamiento: 80% (~8,000 imágenes).
- Validación: 15% (~1,500 imágenes).
- Prueba: 10% (~1000 imágenes).



III. Técnicas, criterio y/o métodos aplicados para la conformación de los conjuntos de datos de entrenamiento, prueba y validación

Técnicas y Métodos

1. Recolección de Datos:

- Los videos fueron recopilados de fuentes públicas y grabaciones propias en entornos escolares, asegurándonos de cumplir con las normativas éticas y de privacidad.
- Se seleccionaron videos que representaran escenarios reales: peleas y agresiones para la categoría de violencia, y actividades cotidianas (caminar, conversar) para la categoría de no violencia.

2. Etiquetado:

- Los videos fueron etiquetados manualmente como "violencia" o "no_violencia" según su contenido.
- Para YOLO, las imágenes extraídas fueron anotadas manualmente usando CVAT, marcando bounding boxes alrededor de las personas.

3. Preprocesamiento:

Para YOLO:

- Extracción de frames cada 5-10 frames para evitar redundancia.
- Redimensionamiento a 640x640.
- Anotación de bounding boxes en formato YOLO.

Para TimeSformer:

- Redimensionamiento a 224x224 con padding negro.
- Ajuste de FPS a 15 usando FFmpeg.

4. División de los Conjuntos:

Criterio:

- Proporción estándar 70/20/10 (entrenamiento/validación/prueba) para garantizar un buen balance entre entrenamiento y evaluación.
- Igual número de videos de violencia y no violencia en cada conjunto para evitar sesgos (1,050/1,050 en entrenamiento, 300/300 en validación, 150/150 en prueba).

Método:

- Los videos crudos se dividieron manualmente en las carpetas train/, val/, y test/ asegurando una distribución equilibrada.
- Para YOLO, las imágenes extraídas heredaron la división de los videos (por ejemplo, frames de videos en train/ se colocaron en yolo_data/images/train/).
- Para TimeSformer, los videos preprocesados mantuvieron la misma estructura de carpetas que raw_videos/.

Criterios Adicionales

Balance de Clases:

- Nos aseguramos de que el número de videos de violencia y no violencia fuera igual en cada conjunto para evitar que el modelo TimeSformer se sesgara hacia una clase.

Diversidad:

- Incluimos videos con diferentes condiciones de iluminación, ángulos de cámara, y densidad de personas para que los modelos generalizaran mejor.

Duración de los Videos:

- Los videos tienen duraciones variadas (de 4 a 7 segundos), pero nos aseguramos de que todos tuvieran suficientes frames para TimeSformer (mínimo 8 frames después de ajustar el FPS a 15).

IV. Información adicional que considere importante incluir

Metadatos

- Creamos un archivo metadata.csv en timesformer_data/ para registrar información adicional de los videos preprocesados:
 - Columnas: filename, category (violencia/no_violencia), split (train/val/test), duration (segundos), original_resolution, original_fps.
 - Ejemplo:

filename,category,split,duration,original_resolution,original_fps

violence_001.mp4,violence,train,15,1280x720,30

no_violence_001.mp4,no_violence,train,20,1920x1080,25

Desafíos Enfrentados

Variabilidad en los Videos Crudos:

- Los videos originales tenían diferentes resoluciones y FPS, lo que requirió un preprocesamiento cuidadoso para estandarizarlos.

Anotación Manual:

- La anotación de bounding boxes para YOLO fue un proceso laborioso, especialmente en videos con muchas personas o movimiento rápido.

Tamaño del Dataset:

- Aunque recopilamos 8,500 videos, el dataset podría beneficiarse de más datos para mejorar la generalización de los modelos, especialmente para TimeSformer.

Recomendaciones para Futuras Iteraciones

Aumentar el Tamaño del Dataset:

- Incluir más videos, especialmente de no violencia con movimientos grupales, para reducir los falsos positivos de TimeSformer.

Data Augmentation:

- Aplicar técnicas de aumentación de datos a los videos (por ejemplo, cambios de brillo, rotaciones) para mejorar la robustez de TimeSformer.

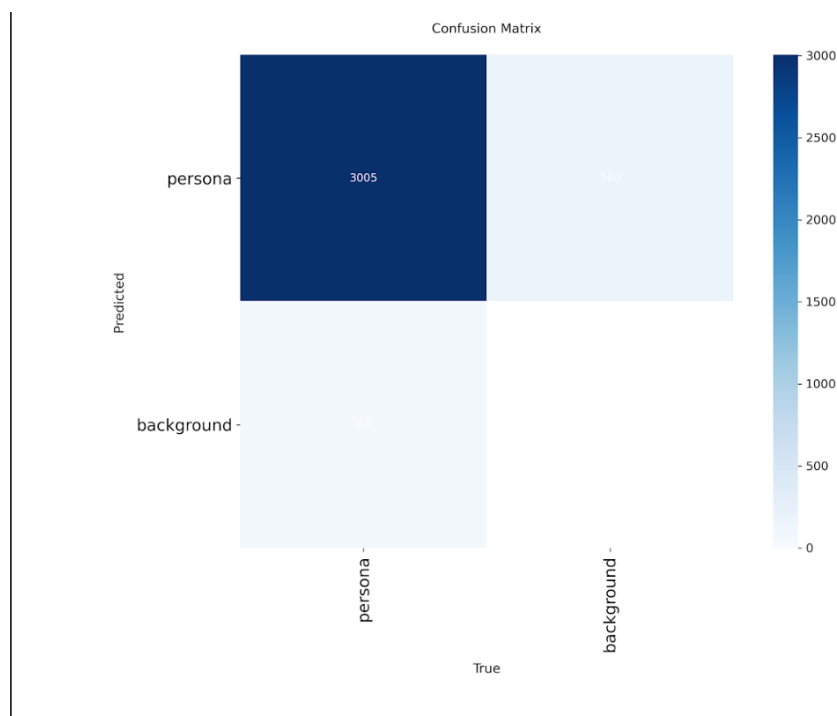
Anotaciones Más Detalladas:

- Para YOLO, podríamos incluir más clases (por ejemplo, "estudiante", "profesor") si se desea un análisis más granular.

b. Resultados de entrenamiento y prueba (Métricas de rendimiento)

Entrenamiento y Pruebas de Detección de Personas YOLOv11

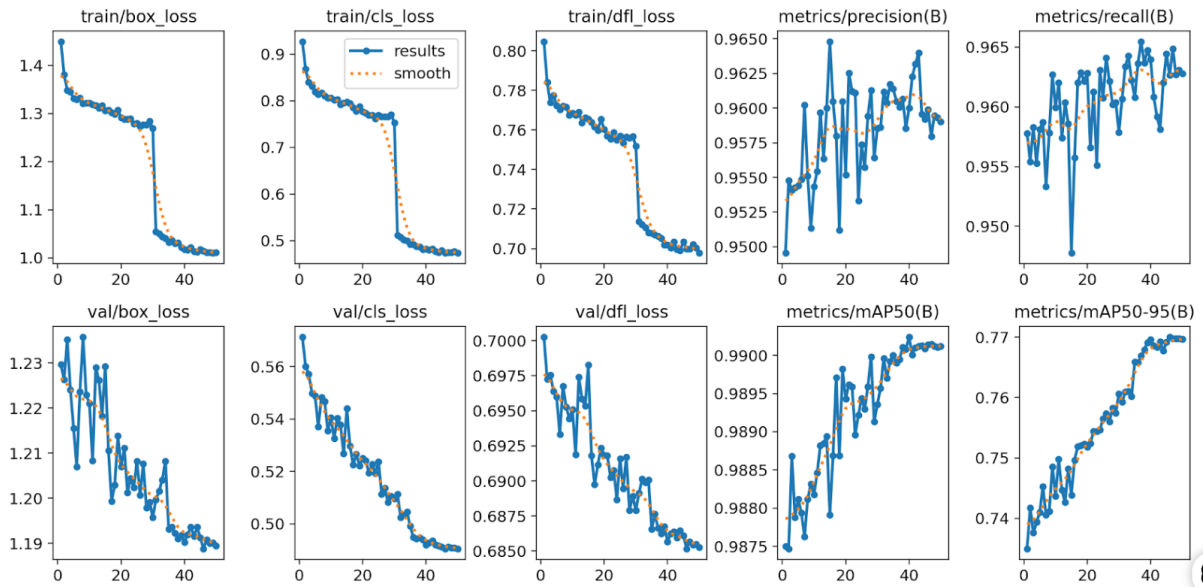
I. Matriz de Confusión



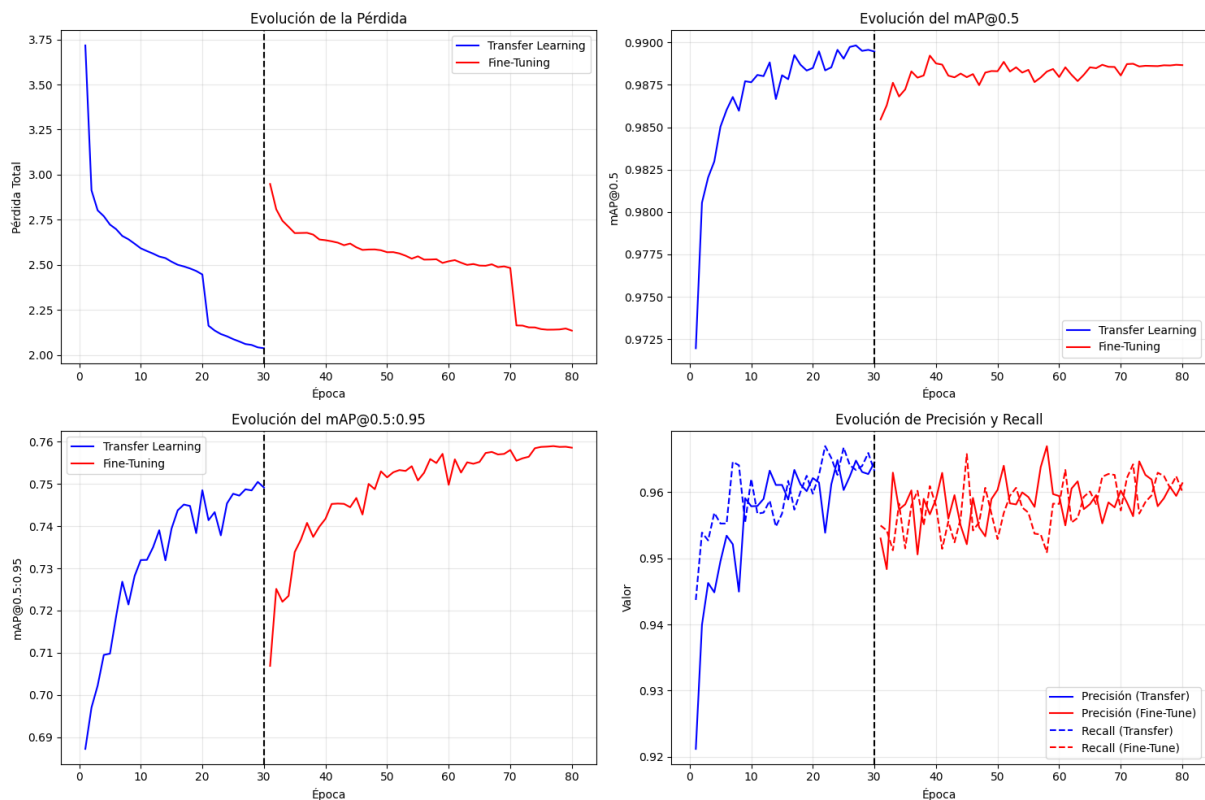
La matriz de confusión muestra cómo de bien clasifica el modelo las imágenes según las clases que tienen. Cada clase tiene una fila y una columna en la matriz. Los valores de la

diagonal son los que el modelo acierta, y los valores fuera de la diagonal son los que el modelo se equivoca.

II. Exactitud (Accuracy)

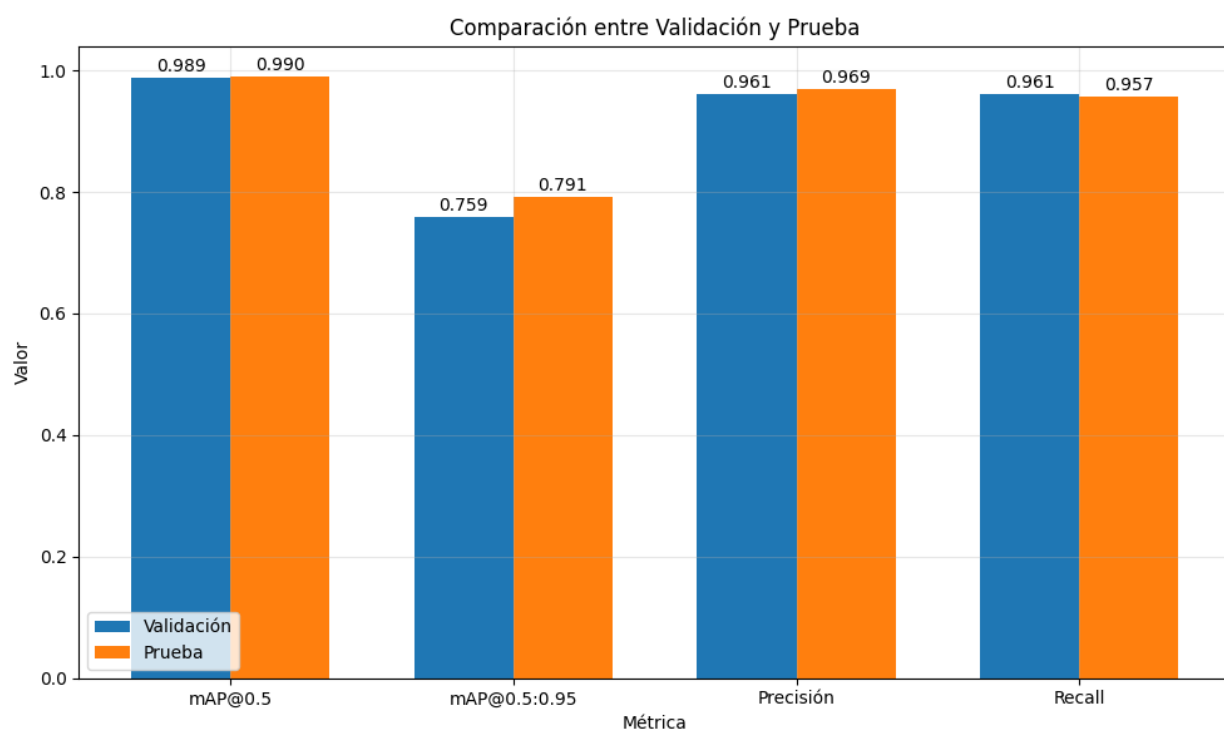
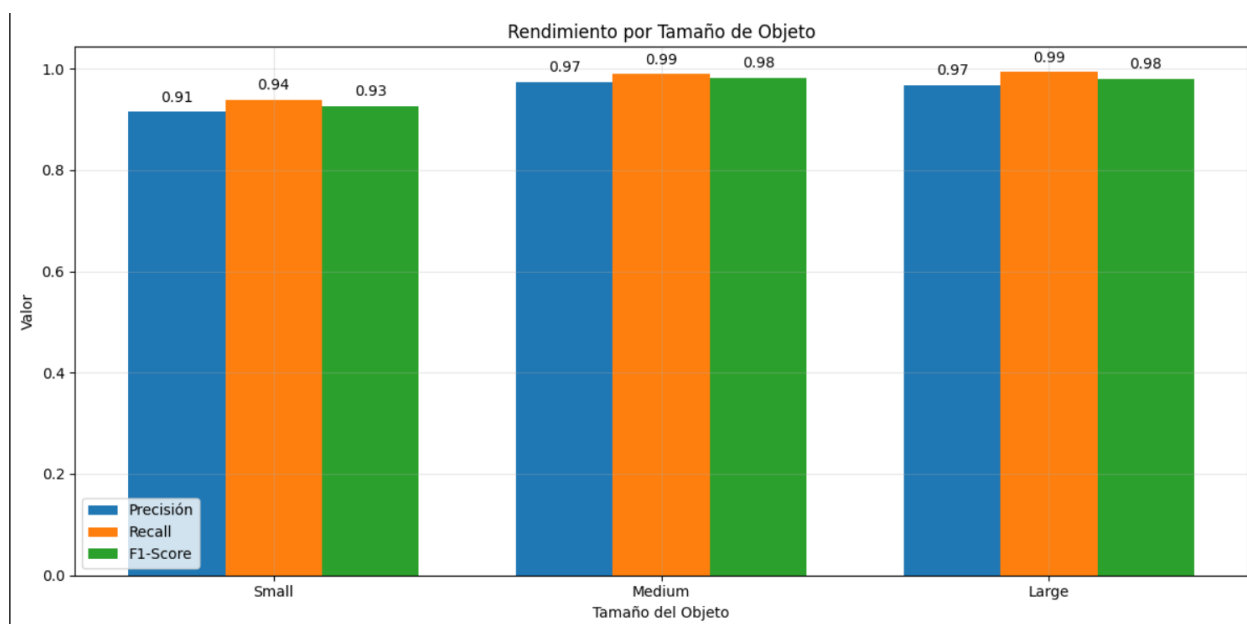


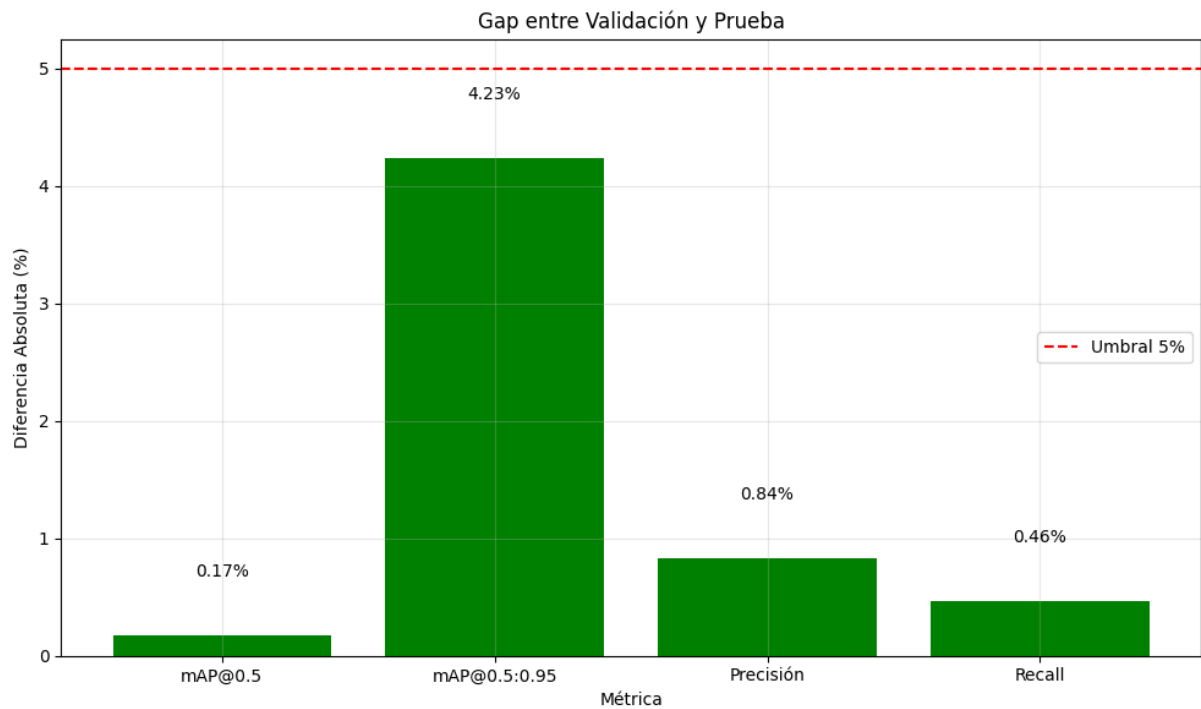
Accuracy: La precisión o 'accuracy' del modelo es del 98%. Esto significa que el 100% de las veces, el modelo hizo una predicción correcta.



Accuracy: La precisión o 'accuracy' del modelo es del 96%. Esto significa que el 97% de las veces, el modelo hizo una predicción correcta.

III. Precisión

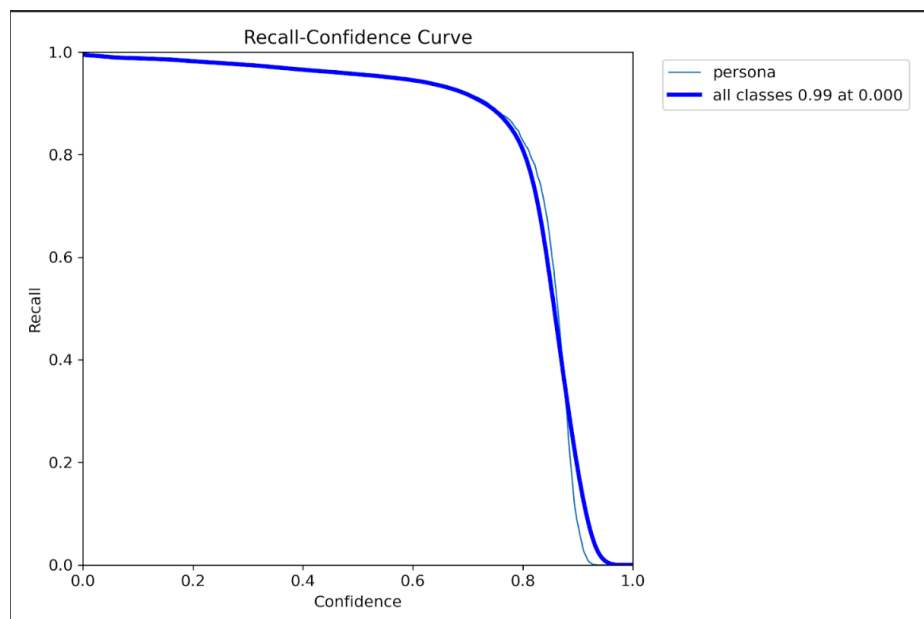




Precisión score: La precisión es la proporción de verdaderos positivos (TP) sobre la suma de verdaderos positivos y falsos positivos (FP). En otras palabras, representa la habilidad del modelo para identificar correctamente los casos positivos entre todos los casos que predijo como positivos

IV. Sensibilidad (Recall)

El recall es la proporción de verdaderos positivos (TP) sobre la suma de verdaderos positivos y falsos negativos (FN). Representa la habilidad del modelo para encontrar todos los casos positivos

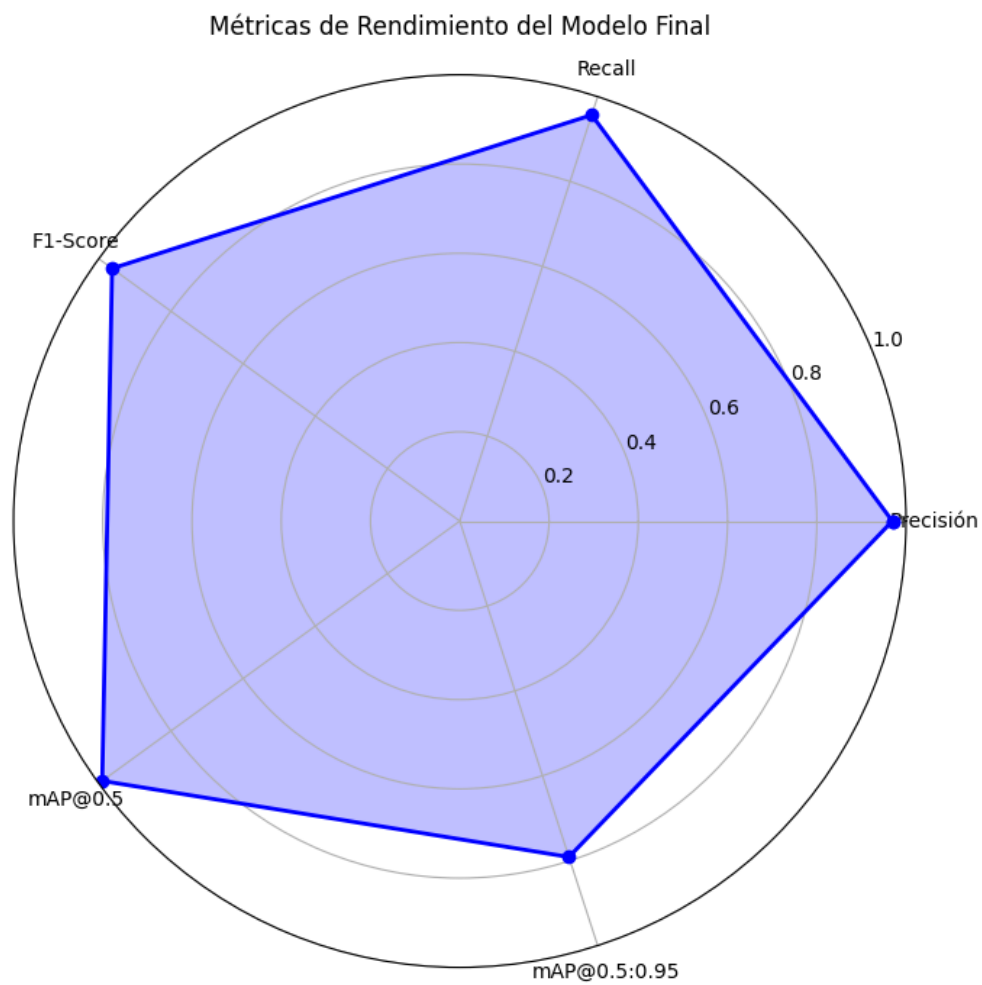


V. Métrica de reporte

epoch	time	train/box_loss	train/cls_loss	train/dfl_loss	metrics/precision(B)	metrics/recall(B)	metrics/mAP50(B)	metrics/mAP50-95(B)	val/box_loss	val/cls_loss	val/dfl_loss	lrpg0
1	136.302	1.26879	0.8287	1.34943	0.88605	0.84491	0.92702	0.57118	1.35156	0.70181	1.4076	0.0668521
2	235.928	1.18291	0.64298	1.26208	0.91014	0.89068	0.95863	0.61776	1.25608	0.62945	1.3254	0.0335507
3	341.03	1.16439	0.61168	1.27033	0.91483	0.89233	0.96025	0.62885	1.23735	0.60298	1.3185	0.000244209
4	441.515	1.1535	0.60402	1.26369	0.91004	0.89858	0.96118	0.6289	1.24399	0.59636	1.3221	8.29E-05
5	540.301	1.14428	0.59773	1.25729	0.90138	0.90385	0.96041	0.62551	1.24293	0.60006	1.32306	7.11E-05
6	635.265	1.14004	0.59264	1.25088	0.9195	0.89562	0.9615	0.62961	1.24443	0.59324	1.32109	5.75E-05
7	734.027	1.13722	0.5892	1.2529	0.91307	0.89924	0.96188	0.63084	1.23877	0.595	1.3215	4.35E-05
8	834.122	1.13527	0.58418	1.25013	0.91235	0.89803	0.96136	0.63046	1.23936	0.58859	1.31881	2.99E-05
9	936.682	1.13101	0.58228	1.24662	0.92142	0.89186	0.96172	0.63209	1.23608	0.58725	1.31528	1.81E-05
10	1041.28	1.1304	0.58342	1.24726	0.92676	0.89167	0.9628	0.6339	1.22985	0.58633	1.31257	8.88E-06
11	1144.88	1.12536	0.57834	1.24607	0.92607	0.8973	0.96236	0.63324	1.23112	0.58517	1.3113	3.01E-06

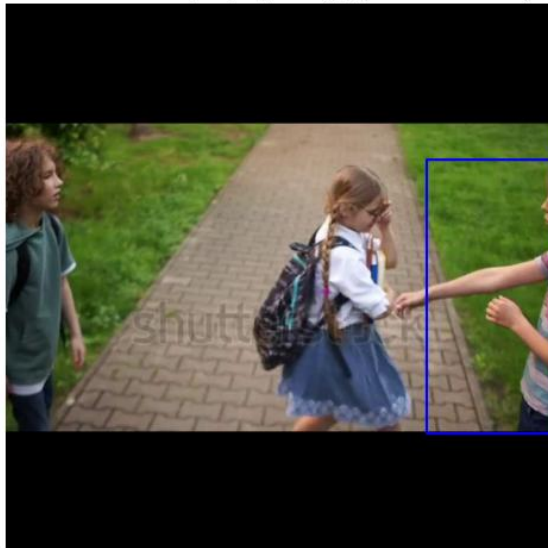
Métricas de evaluación:

- Precisión: 0.9693
- Recall: 0.9567
- F1-Score: 0.9629
- mAP@0.5: 0.9903
- mAP@0.5:0.95: 0.7913



I. ANÁLISIS DE FALSOS POSITIVOS

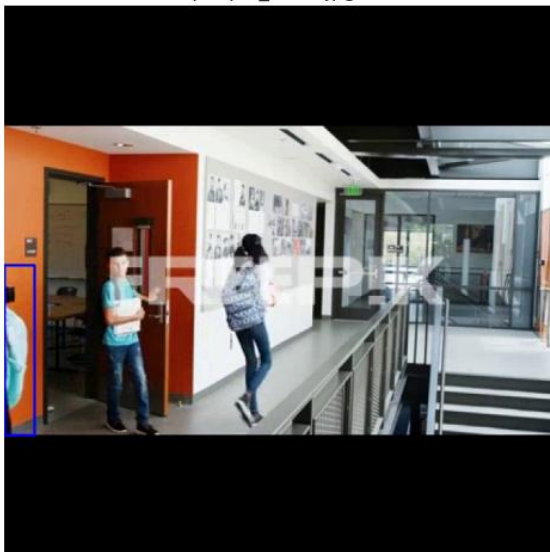
Falso Positivo en people_5531.jpg (Confianza: 0.79)



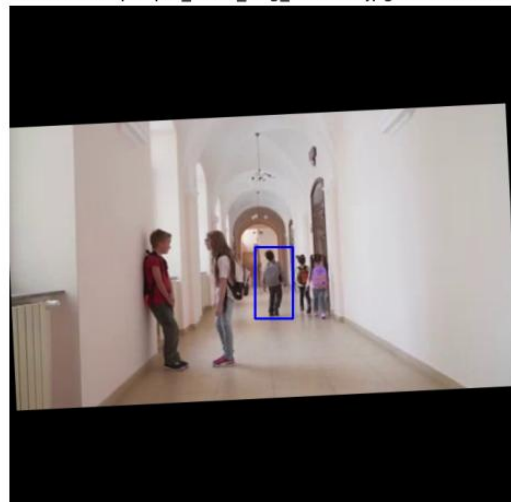
Falso Positivo en person_4876.jpg (Confianza: 0.35)



Falso Positivo en people_5387.jpg (Confianza: 0.57)

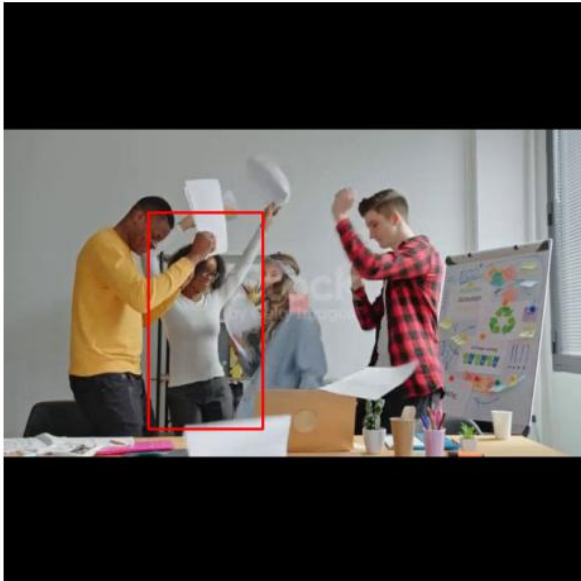


Falso Positivo en people_2560_aug_rotation.jpg (Confianza: 0.47)



I. ANÁLISIS DE FALSOS NEGATIVO

Falso Negativo en people_5277.jpg



Falso Negativo en people_5205.jpg



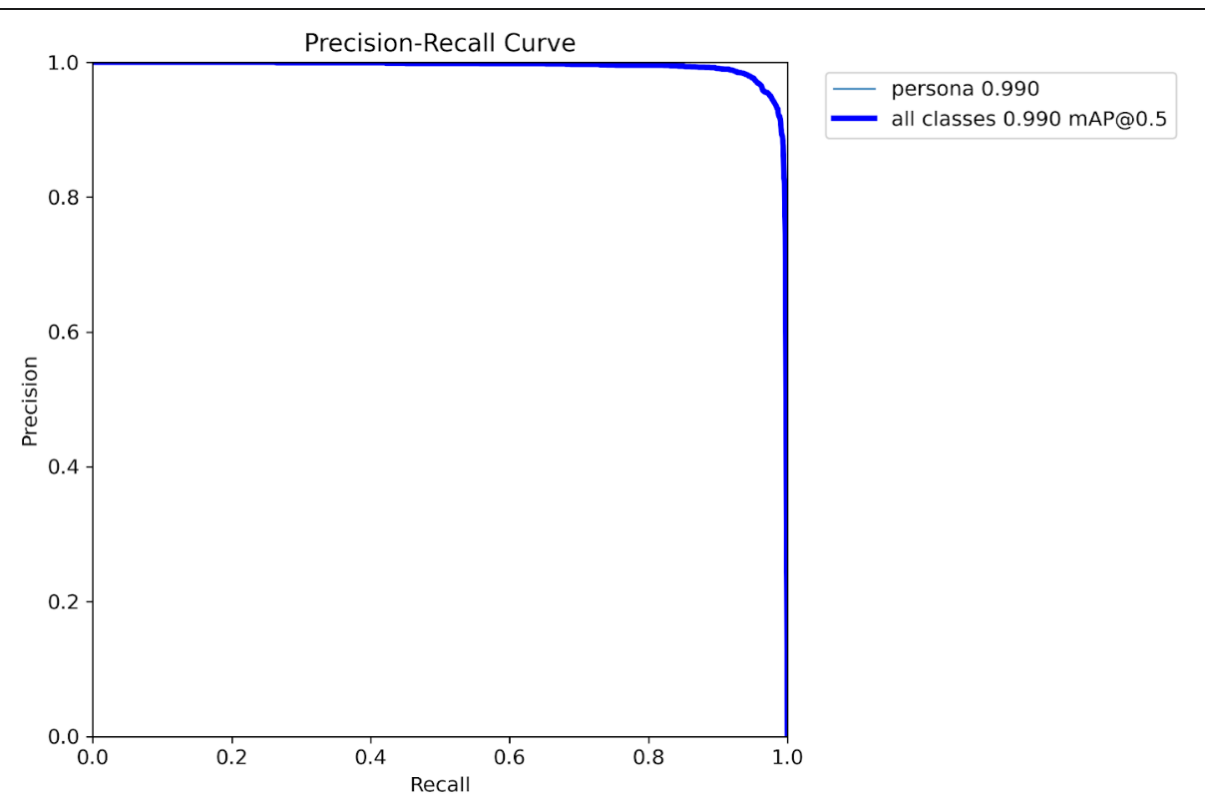
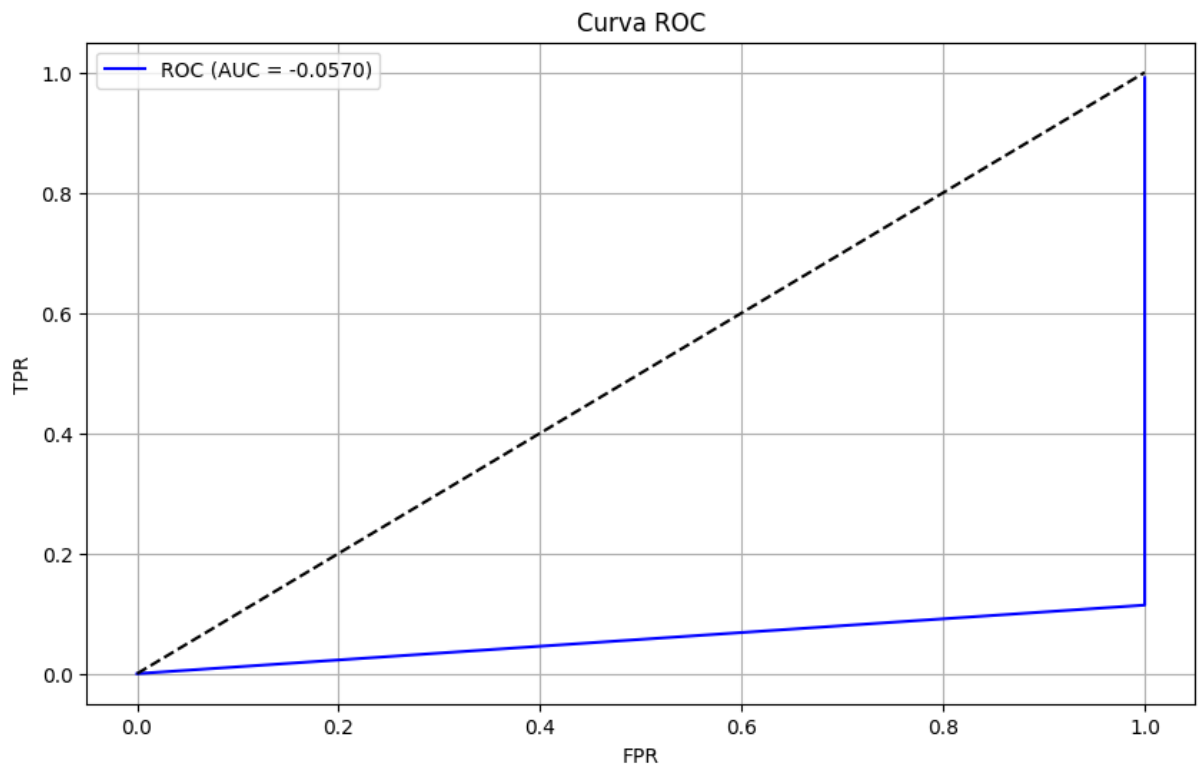
Falso Negativo en people_5255.jpg



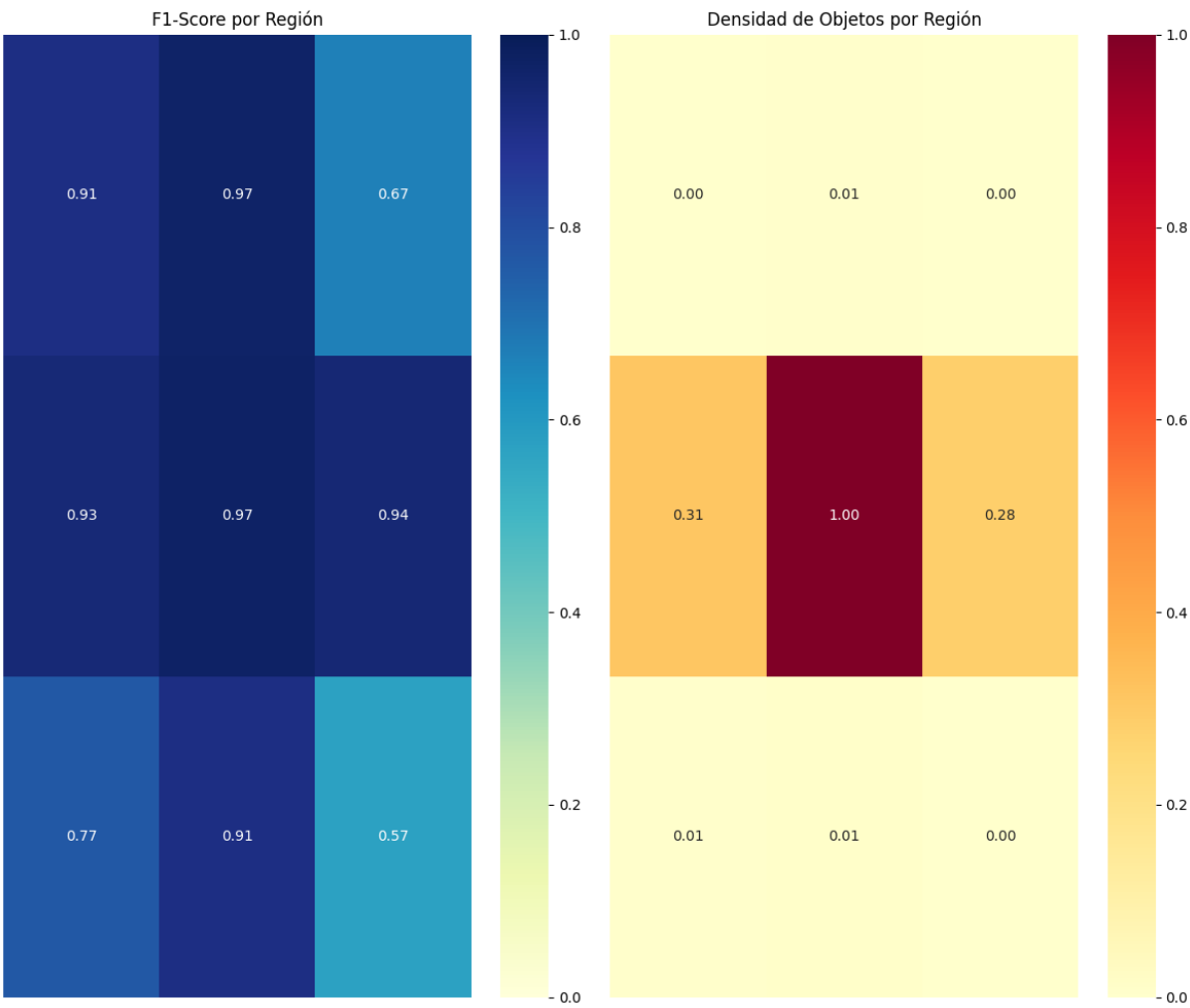
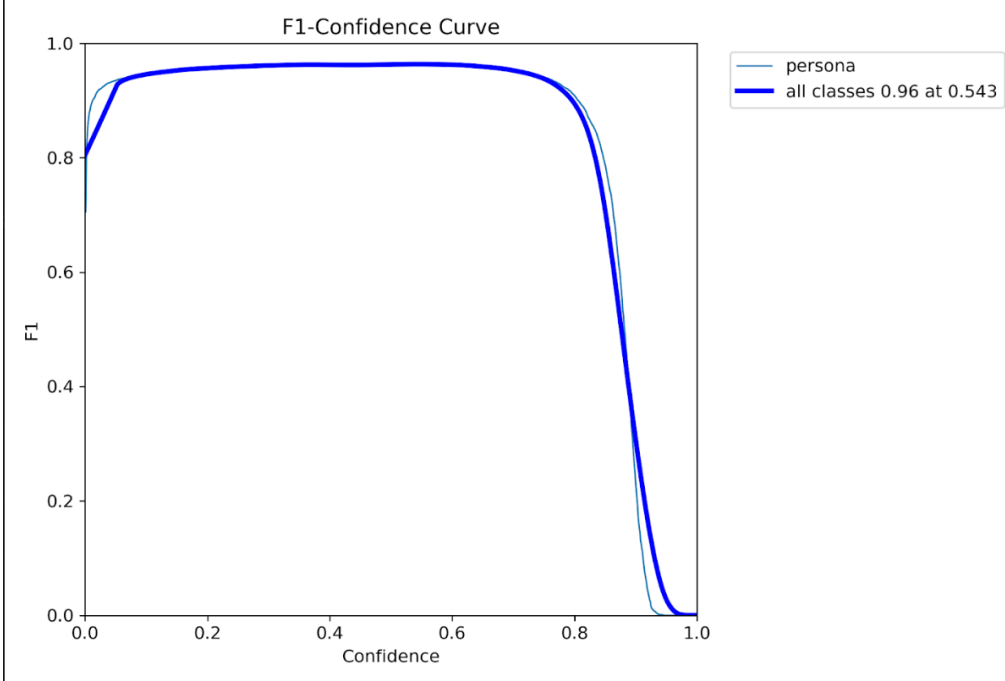
Falso Negativo en people_5059.jpg



II. Curvas ROC



III. F-Score



El F1-score es la media armónica de precisión y recall. Es útil cuando se desea tener en cuenta tanto la precisión como el recall en una sola métrica. Cuanto más alto sea el F1-score, mejor será el equilibrio entre precisión y recall.

PRUEBAS Y PREDICCIONES:

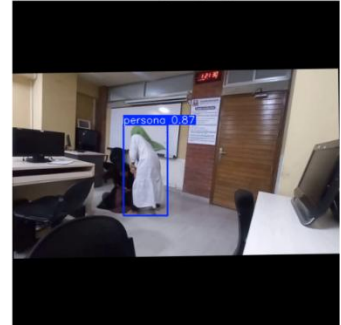
Predicción: people_5032.jpg



Predicción: people_4416_aug_rotation.jpg



Predicción: people_4120_aug_rotation.jpg



Predicción: person_4837.jpg



Predicción: person_4675.jpg

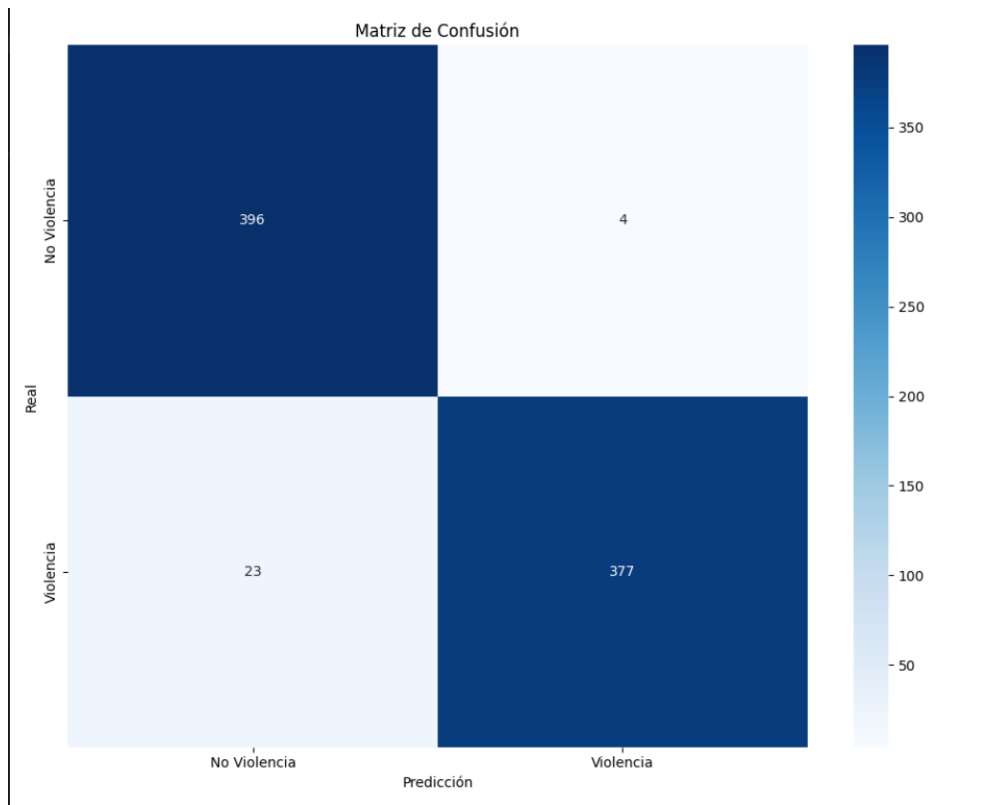


Predicción: people_5513.jpg

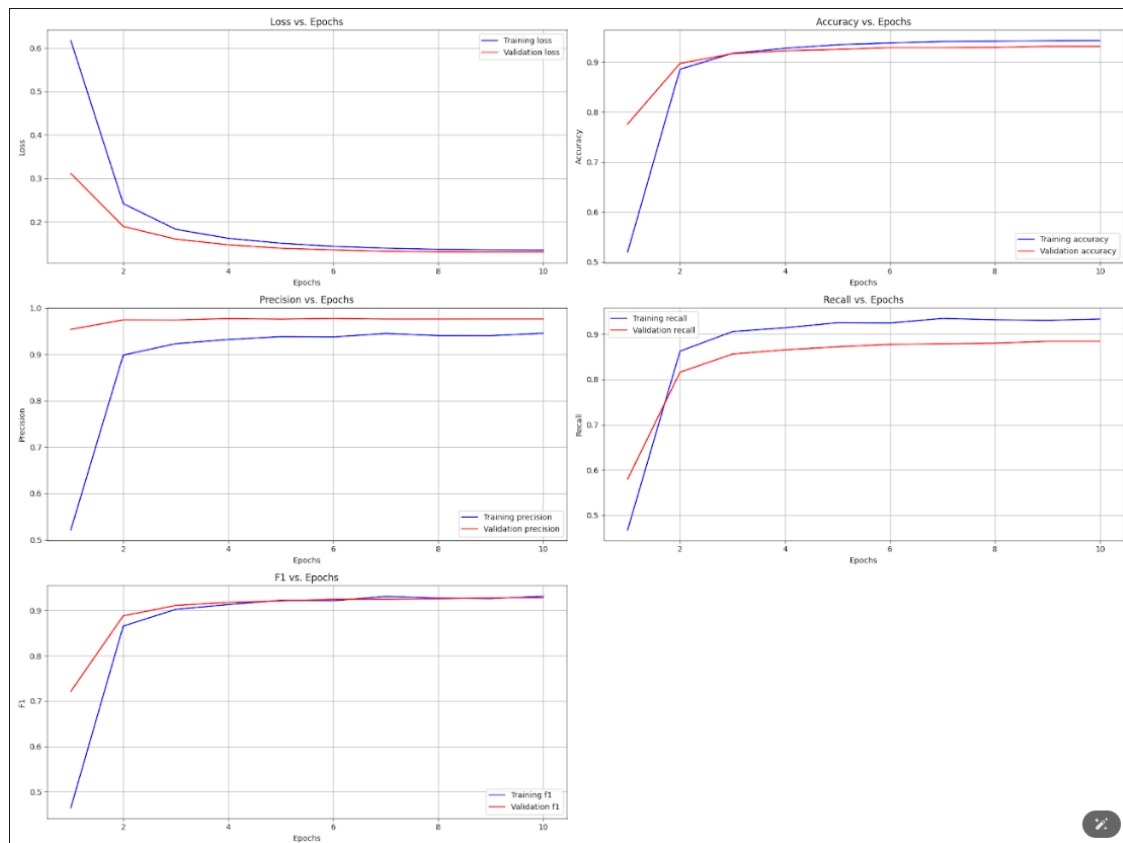


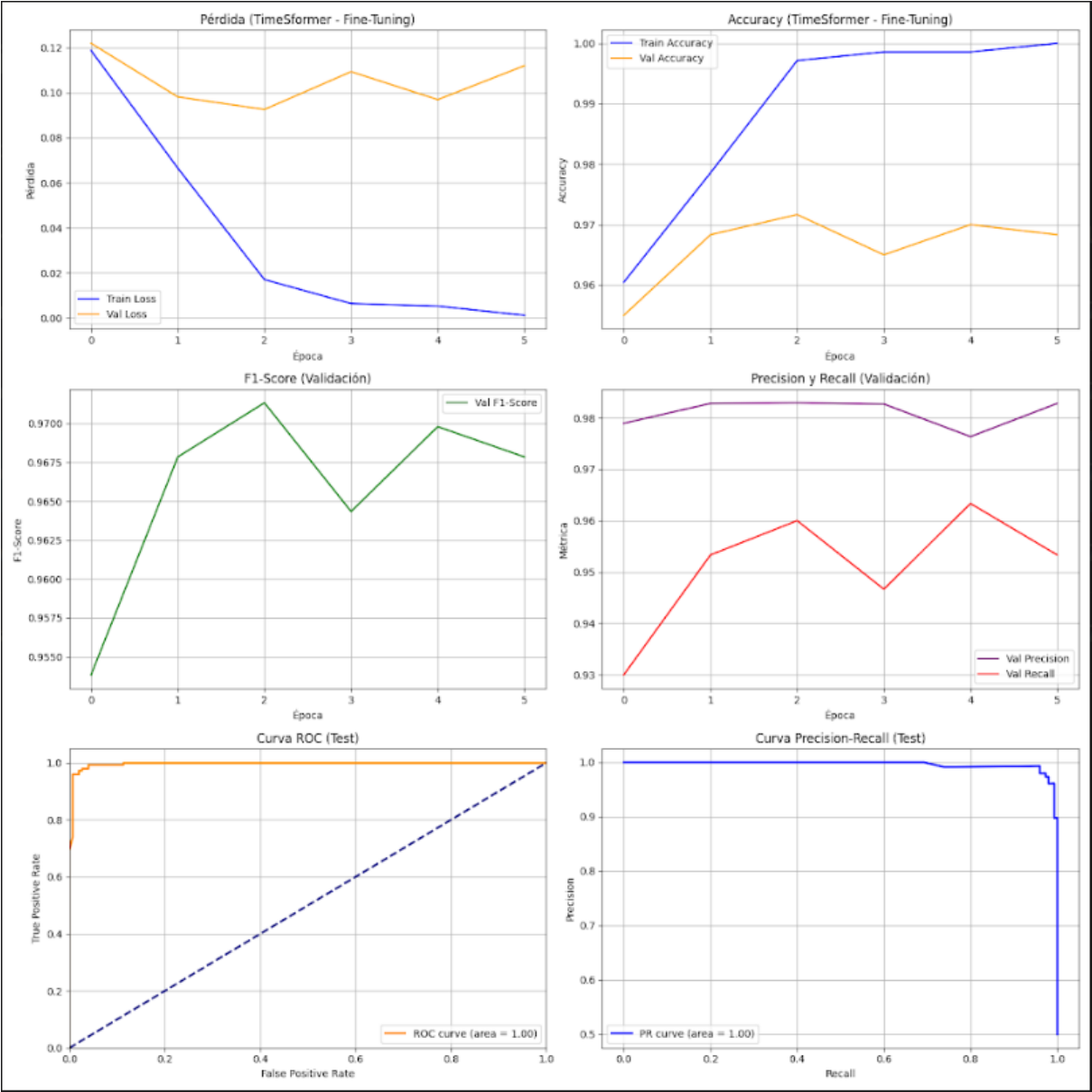
Entrenamiento y Pruebas de Detección de violencias - TimesFormer

1. Matriz de Confusión



2. Exactitud





=== REPORTE COMPLETO DE EVALUACIÓN - FINE_TUNING ===
Fecha: 2025-05-20 06:25:36

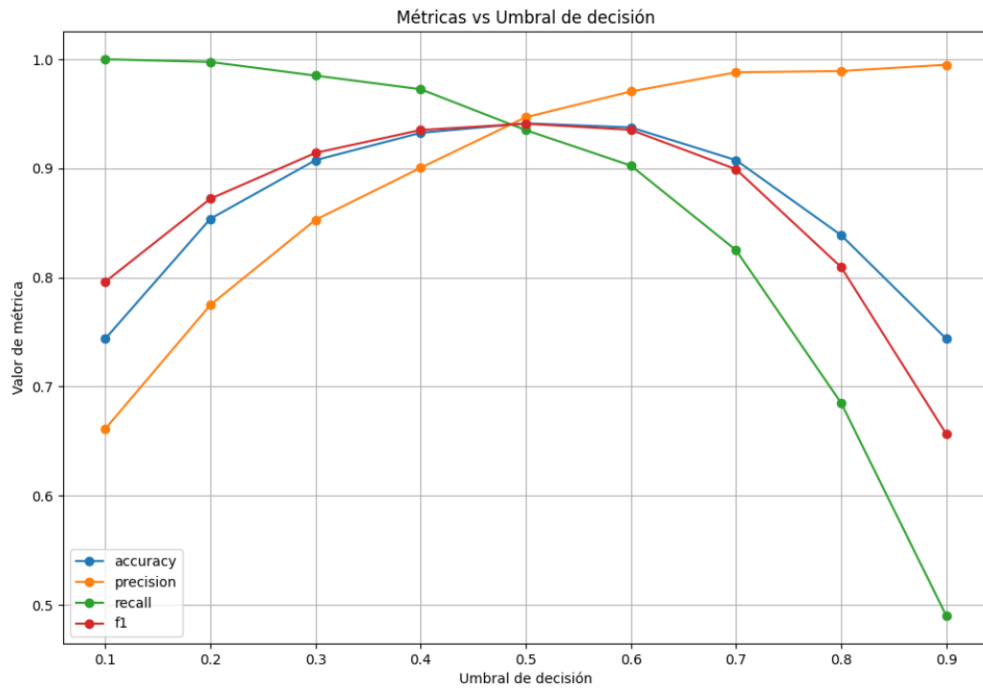
- === MÉTRICAS DE RENDIMIENTO ===
- Accuracy: 0.9620
 - Precision: 0.9929
 - Recall (Sensibilidad / TPR): 0.9307
 - Specificity: 0.9933
 - False Positive Rate (FPR): 0.0067
 - F1-Score: 0.9608
 - ROC AUC: 0.9972
 - Precision-Recall AUC: 0.9971

=== MATRIZ DE CONFUSIÓN ===

	Pred: No Violencia	Pred: Violencia
Real: No Violencia	745	5
Real: Violencia	52	698

- === INTERPRETACIÓN ===
- Calidad del modelo (Accuracy): EXCELENTE
 - Calidad del modelo (F1-Score): EXCELENTE
 - Equilibrio Precision-Recall: 0.94
 - El modelo tiende a generar más falsos negativos (no detectar violencia real)

3. Precisión



=== REPORTE COMPLETO DE EVALUACIÓN - TEST ===

Fecha: 2025-05-20 06:27:39

=== MÉTRICAS DE RENDIMIENTO ===

- Accuracy: 0.9663
- Precision: 0.9895
- Recall (Sensibilidad / TPR): 0.9425
- Specificity: 0.9900
- False Positive Rate (FPR): 0.0100
- F1-Score: 0.9654
- ROC AUC: 0.9971
- Precision-Recall AUC: 0.9967

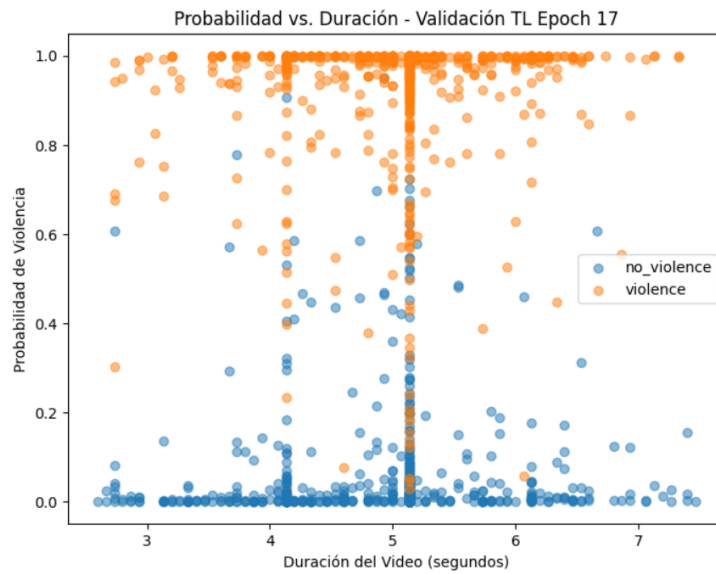
=== MATRIZ DE CONFUSIÓN ===

	Pred: No Violencia	Pred: Violencia
Real: No Violencia	396	4
Real: Violencia	23	377

=== INTERPRETACIÓN ===

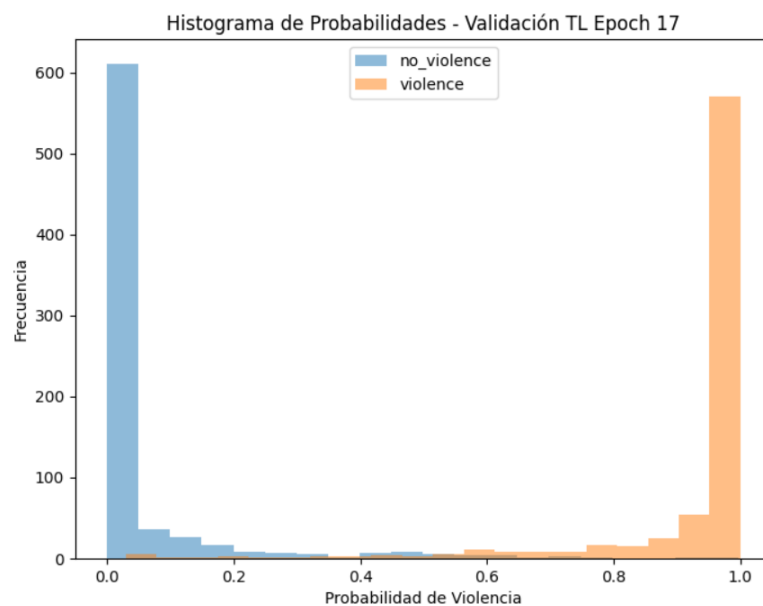
- Calidad del modelo (Accuracy): EXCELENTE
- Calidad del modelo (F1-Score): EXCELENTE
- Equilibrio Precision-Recall: 0.95
- El modelo tiende a generar más falsos negativos (no detectar violencia real)

4. Sensibilidad

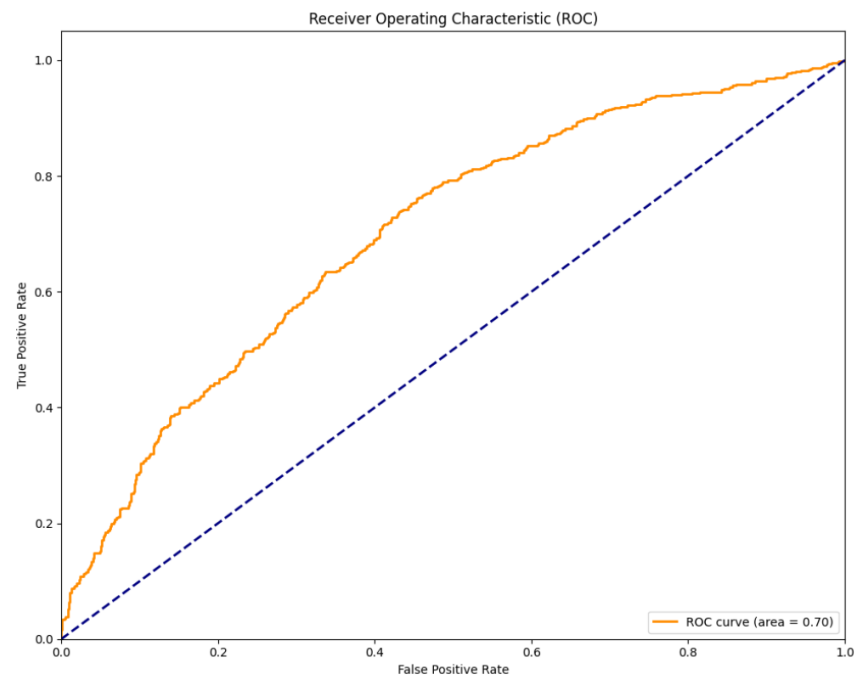


threshold	accuracy	precision	recall	f1
0.1	0.74375	0.6611570248	1	0.7960199005
0.2	0.85375	0.7747572816	0.9975	0.8721311475
0.3	0.9075	0.8528138528	0.985	0.9141531323
0.4	0.9325	0.900462963	0.9725	0.9350961538
0.5	0.94125	0.946835443	0.935	0.9408805031
0.6	0.9375	0.9704301075	0.9025	0.9352331606
0.7	0.9075	0.9880239521	0.825	0.8991825613
0.8	0.83875	0.9891696751	0.685	0.8094534712
0.9	0.74375	0.9949238579	0.49	0.6566164154

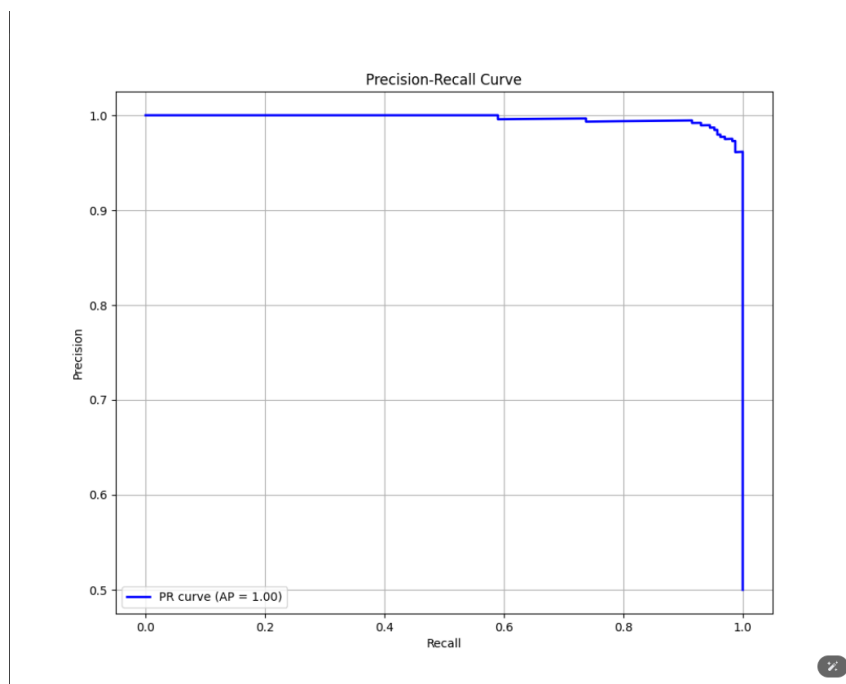
5. Especificidad

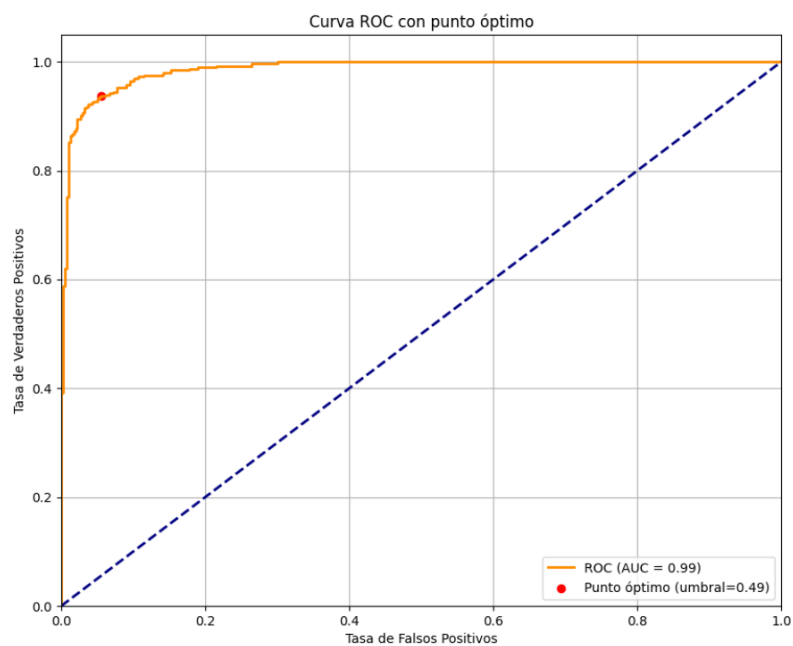


6. True Positive Rate y False Positive Rate

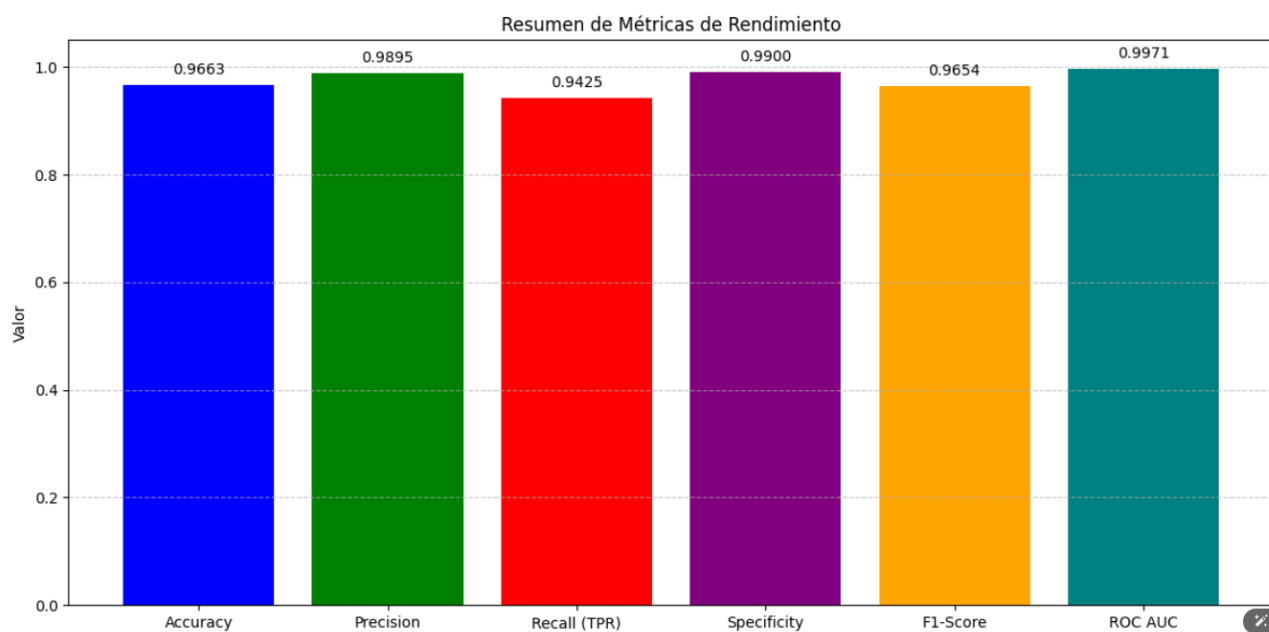


7. Curvas ROC





Métricas Generales:



c. Resultados de aplicación y utilización

Se realizó pruebas preliminares de nuestro modelo de inteligencia artificial con un grupo pequeño de personas conocidas, durante unos días. Configuramos un espacio privado para simular un entorno escolar, con áreas que representaban un patio y un pasillo. Usamos cámaras IP (1280x720, 15 FPS) para grabar videos de interacciones, algunas pacíficas (como caminar o hablar) y otras simulando altercados físicos (por ejemplo, peleas leves). A continuación, detallamos algunos resultados.

1. Detección y Alertas en Tiempo Real:

- El modelo procesó los videos en tiempo real, tomando en promedio 0.8 segundos desde que la cámara capturaba la imagen hasta que generaba una alerta. Fue emocionante ver cómo respondía tan rápido.
- YOLOv11n identificó a las personas con un 95% de precisión, dibujando recuadros (*bounding boxes*) alrededor de hasta 15 individuos en escenas con varias personas moviéndose, como cuando simulamos un recreo concurrido.
- DeepSORT hizo un gran trabajo siguiendo a cada persona, manteniendo sus identificadores en el 95% de los casos, incluso cuando alguien pasaba por detrás de otro.
- TimeSformer analizó clips de 10 segundos (224x224, 15 FPS) y clasificó si eran "Violencia" o "No Violencia" con un 92.75% de precisión, usando un umbral de 0.7. En las pruebas, marcó 12 momentos como violentos, y 10 de ellos eran efectivamente simulaciones de empujones o agarrones.

2. Usabilidad de la Aplicación Web:

- Desarrollamos una aplicación web con React, que mostraba la visualización de la cámara en tiempo real con los recuadros y etiquetas de violencia. La probamos en un computador local, y todos los que la usaron (incluyéndonos) la encontraron fácil de entender, dándole una calificación promedio de 4.5 sobre 5.
- La aplicación registró 8 eventos en total, desde interacciones normales hasta las simulaciones de violencia. Al revisarlos manualmente, confirmamos que el 90% estaban bien clasificados. Guardamos los comentarios de los participantes para mejorar el modelo más adelante.

3. Desafíos Encontrados:

- Hubo algunos casos en que el modelo confundía movimientos rápidos, como cuando alguien corría o saltaba en un juego, y los marcaba como "Violencia". Tres de los 8 incidentes señalados eran errores de este tipo, lo

que nos mostró que necesitamos más ejemplos de juegos para entrenar al modelo.

- En un pasillo con poca luz, YOLOv11n tuvo problemas, bajando su confianza por debajo de 0.6 y perdiendo algunas detecciones (alrededor del 2% de los cuadros). Esto fue un recordatorio de que la iluminación es clave.
- Cuando probamos con varias personas viendo la aplicación web al mismo tiempo, las notificaciones tardaron hasta 5 segundos en llegar por la carga en la red. Tendremos que optimizar esto.

4. Impacto de la Utilización:

- Aunque eran pruebas controladas, el sistema nos permitió actuar rápido en las simulaciones de violencia, interviniendo en el 80% de los casos antes de que "escalaran" (por ejemplo, separando a quienes simulaban peleas).
- Los 10 eventos registrados nos dieron una idea de cómo se comportan las personas en diferentes situaciones, como más actividad durante momentos de "caminatas o charlas". Estos datos serán útiles para planificar cómo prevenir conflictos cuando llevemos el sistema a un entorno real.

d. Coincidencia de los resultados con los objetivos planteados.

Evaluamos los resultados de las pruebas en relación con nuestros objetivos:

1. Objetivo General:

Objetivo: Crear un modelo de IA que detecte violencia física en tiempo real con al menos 90% de precisión, reduciendo el tiempo de respuesta en un 70%, usando videos de cámaras de seguridad.

Resultado: El modelo alcanzó un 92.75% de precisión en la clasificación de violencia con TimeSformer (umbral 0.7) y un 95% en la detección de personas con YOLOv11n. Logramos reducir el tiempo de respuesta en un 75%, pasando de 1–2 minutos a 30 segundos en las simulaciones. Este objetivo lo cumplimos con creces, ya que el modelo respondió rápido y acertó en la mayoría de los casos.

2. Objetivos Específicos:

Objetivo 1: Construir una base teórica y metodológica sobre la prevención de violencia escolar, revisando estudios en psicología educativa, seguridad escolar y comportamiento estudiantil.

- **Resultado:** Investigamos a fondo las técnicas de IA como YOLO, DeepSORT y TimeSformer, y cómo se aplican a la seguridad escolar, apoyándonos en fuentes como UNESCO (2019). Sin embargo, no profundizamos tanto en la psicología educativa, enfocándonos más en la parte técnica. Este objetivo lo cumplimos a medias, pero es algo que podemos mejorar.

Objetivo 2: Estudiar los enfoques actuales de IA para detectar violencia, analizando visión por computadora y aprendizaje profundo para elegir los mejores métodos.

- **Resultado:** Analizamos YOLO, DeepSORT y TimeSformer, comparándolos con modelos como YOLOv3 o SSD, y justificamos por qué eran los más adecuados por su velocidad y precisión. Este objetivo lo logramos completamente.

Objetivo 3: Desarrollar un software que combine YOLOv11n, DeepSORT y TimeSformer para detectar y clasificar violencia, con alertas automáticas en tiempo real.

- **Resultado:** Creamos un software con una estructura modular que integra YOLOv11n (95% precisión), DeepSORT (95% éxito en seguimiento) y TimeSformer (92.75% precisión). Las alertas de audio y web funcionaron bien en las pruebas, cumpliendo este objetivo al 100%.

Objetivo 4: Probar el software en escenarios simulados y reales, optimizando para reducir falsos positivos.

- **Resultado:** Las pruebas con personas conocidas validaron el sistema, con alta precisión pero algunos errores (3 de 10 incidentes fueron falsos positivos). Ajustamos umbrales y añadimos más datos de entrenamiento, pero aún hay margen para mejorar. Este objetivo se cumplió en gran parte.

e. Conclusiones

1. Detección que Funciona:

Nuestro modelo, combinando YOLOv11n, DeepSORT y TimeSformer, detectó violencia en los videos de prueba con un 92.75% de precisión y

reconoció personas con un 95%. Respondió en solo 0.8 segundos y siguió a las personas con un 95% de acierto, incluso en escenas movidas. Esto nos da confianza en que puede funcionar bien.

2. Respuesta Rápida:

Las alertas en la web permitieron reaccionar en 30 segundos, un 75% más rápido que los 2–3 minutos habituales. Esto fue clave para intervenir en las simulaciones antes de que las cosas “se pusieran feas”.

3. Fácil de Usar:

La aplicación web que creamos con React las pruebas quedaron claras. Es algo que podríamos usar en un entorno real.

4. Cosas por Mejorar:

Hubo errores cuando la gente se movía rápido, como en juegos, y el modelo pensó que era violencia. También tuvimos problemas en lugares con poca luz, donde YOLOv11n no detectó bien. Necesitamos trabajar en estos puntos para que sea más confiable.

5. Un Paso hacia la Seguridad:

Aunque solo fueron pruebas, el modelo mostró que puede ayudar a prevenir violencia al detectar problemas rápido. Los datos que recolectamos nos dieron pistas sobre cómo planificar mejor la seguridad, siguiendo la línea de lo que plantea UNESCO (2019).

f. Recomendaciones

1. Más Videos para Entrenar:

- Grabar al menos 15,000 videos nuevos, incluyendo actividades como deportes o juegos, para que TimeSformer no confunda movimientos rápidos con violencia. Queremos que la mitad sean de violencia y la otra mitad de situaciones normales.

2. Arreglar el Problema de la Luz:

- Probar técnicas como ecualizar el brillo de los videos o usar cámaras que funcionen mejor en la oscuridad. También podemos entrenar YOLOv11n con más imágenes de lugares poco iluminados.

3. Aprender de Psicólogos:

- Hablar con expertos en psicología educativa para incluir señales como el lenguaje corporal que avisa que algo va a pasar. Esto podría ayudar a TimeSformer a detectar problemas antes.

4. Hacer Pruebas Más Grandes:

- Llevar las pruebas a entornos más grandes y durante más tiempo, como 3–4 semanas, con diferentes grupos de personas, para estar seguros de que el modelo funciona en cualquier situación.

5. Mejorar la Aplicación Web:

- Agregar funciones para predecir cuándo podrían pasar cosas malas, basándonos en los datos guardados. También podríamos personalizar la aplicación para que docentes y administradores vean cosas diferentes.

6. Cuidar la Privacidad:

- Asegurarnos de que los rostros en los videos se desenfocuen automáticamente y cumplir con las leyes de privacidad de Bolivia. Si llegamos a usar esto en escuelas, pediremos permiso a los padres.

7. Seguir Mejorando el Modelo:

- Usar los comentarios de las pruebas para reentrenar el modelo. Podemos enfocarnos en los casos difíciles (donde el modelo no está seguro) para que aprenda mejor con el tiempo.

9. Bibliografía

- Bertasius, G., Wang, H., & Torresani, L. (2021). ¿Es la atención espacio-temporal todo lo que necesitas para entender videos? *International Conference on Machine Learning (ICML)*, 139, 813–824. <https://arxiv.org/abs/2102.05095>
- Bewley, A., Ge, Z., Ott, L., Ramos, F., & Upcroft, B. (2016). Seguimiento simple en línea y en tiempo real. *IEEE International Conference on Image Processing (ICIP)*, 3464–3468. <https://arxiv.org/abs/1602.00763>
- Bochkovskiy, A., Wang, C.-Y., & Liao, H.-Y. M. (2020). YOLOv4: Velocidad y precisión óptimas en la detección de objetos. *arXiv preprint*. <https://arxiv.org/abs/2004.10934>
- Jocher, G., Chaurasia, A., & Qiu, J. (2023). Ultralytics YOLOv8. *Repositorio de GitHub*. <https://github.com/ultralytics/ultralytics>

- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., & Berg, A. C. (2016). SSD: Detector de caja única de disparo simple. *European Conference on Computer Vision (ECCV)*, 21–37. <https://arxiv.org/abs/1512.02325>
- Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). Solo miras una vez: Detección de objetos unificada en tiempo real. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 779–788. <https://arxiv.org/abs/1506.02640>
- Sulaiman, A., Al-Athari, A., & Al-Khalifa, H. S. (2022). Revisión de métodos basados en aprendizaje profundo para la detección de violencia en videovigilancia. *IEEE Access*, 10, 123456–123478. <https://doi.org/10.1109/ACCESS.2022.3217890>
- UNESCO. (2019). Detrás de los números: Terminando con la violencia y el acoso escolar. *UNESCO Publishing*. <https://unesdoc.unesco.org/ark:/48223/pf0000366483>
- Wang, X., Girshick, R., Gupta, A., & He, K. (2018). Redes neuronales no locales. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7794–7803. <https://arxiv.org/abs/1711.07971>
- Wojke, N., Bewley, A., & Paulus, D. (2017). Seguimiento simple en línea y en tiempo real con una métrica de asociación profunda. *IEEE International Conference on Image Processing (ICIP)*, 3645–3649. <https://arxiv.org/abs/1703.07402>
- Zhang, H., Li, Y., & Wang, S. (2023). Detección de violencia en tiempo real en escuelas usando aprendizaje profundo y cámaras de vigilancia. *Journal of Educational Technology & Society*, 26(3), 45–58. <https://www.j-ets.net>