

Evaluación de la Técnica del Bulgarian Split Squat Basada en Aprendizaje Profundo Usando MediaPipe y Redes Recurrentes Bidireccionales

Juan Jose Núñez, Juan Jose Castro
Universidad San Buenaventura
Cali, Colombia

Resumen—La rehabilitación física y la técnica adecuada en el ejercicio son cruciales para prevenir lesiones musculoesqueléticas y optimizar el rendimiento atlético. Los métodos tradicionales de evaluación dependen en gran medida de la supervisión experta, lo cual es costoso en tiempo, subjetivo y no escalable para entrenamiento domiciliario. Este artículo presenta un sistema automatizado de visión por computadora para la evaluación en tiempo real de la ejecución del Bulgarian Split Squat utilizando estimación de pose con MediaPipe y redes Bidireccionales Recurrentes con Unidades Gateadas (BiGRU). Nuestro enfoque extrae 33 puntos anatómicos por fotograma y calcula características biomecánicas incluyendo ángulos articulares, rango de movimiento (ROM) y suavidad del movimiento. El clasificador BiGRU mejorado con mecanismo de atención logra un macro-F1 de 66.38 % y una precisión de 98.37 % en un conjunto de datos de 820 repeticiones extraídas de 74,171 fotogramas, identificando correctamente cuatro patrones de error: ejecución correcta (E0), inclinación excesiva del tronco (E1), valgo de rodilla (E2) y profundidad insuficiente (E3). El hallazgo metodológico más significativo es el impacto crítico del split estratificado: garantizar representación balanceada de clases en entrenamiento, validación y prueba mejoró la precisión de 20.9 % a 98.37 % (+370 % relativo), demostrando que la división de datos es tan importante como la arquitectura del modelo en tareas con desbalance extremo. El sistema demuestra un rendimiento excelente en la clase de error dominante E1 (F1=99.13 %, precisión=98.28 %) manteniendo capacidades de inferencia en tiempo real (<50ms por fotograma), proporcionando retroalimentación inmediata para rehabilitación domiciliaria.

Index Terms—Visión por computadora, estimación de pose, evaluación de ejercicio, BiGRU, MediaPipe, rehabilitación, análisis de calidad de movimiento

I. INTRODUCCIÓN

El ejercicio físico es fundamental para mantener la salud musculoesquelética, prevenir enfermedades crónicas y mejorar la calidad de vida [1], [2]. El Bulgarian Split Squat es un ejercicio unilateral de tren inferior ampliamente prescrito en protocolos de rehabilitación y programas de entrenamiento de fuerza debido a su efectividad en desarrollar la fuerza del cuádriceps, mejorar el equilibrio y abordar asimetrías bilaterales [3]. Sin embargo, una ejecución incorrecta—como inclinación excesiva del tronco hacia adelante, valgo de rodilla (colapso hacia adentro) o profundidad insuficiente—puede conducir a patrones de movimiento compensatorios, eficacia reducida del entrenamiento y mayor riesgo de lesión [4].

La evaluación tradicional de la técnica del ejercicio se basa en la observación manual por clínicos o entrenadores

capacitados, un proceso que es subjetivo, intensivo en tiempo y requiere experiencia especializada. Este enfoque presenta barreras significativas para individuos comprometidos en rehabilitación domiciliaria o entrenamiento remoto, donde la supervisión profesional es limitada o no disponible [5], [6]. Los avances recientes en visión por computadora y aprendizaje profundo ofrecen soluciones prometedoras para automatizar la evaluación de la calidad del movimiento, proporcionando herramientas objetivas, escalables y accesibles para retroalimentación en tiempo real.

I-A. Brecha de Investigación y Limitaciones del Trabajo Previo

Aunque numerosos estudios han explorado la evaluación automatizada de ejercicios usando cámaras RGB y sensores de profundidad [1], [2], persisten varias limitaciones críticas:

- **Deficiencias en modelado temporal:** Muchos enfoques se basan en clasificación fotograma por fotograma o reglas artesanales basadas en umbrales angulares [7], [8], fallando en capturar la dinámica temporal y dependencias secuenciales inherentes al movimiento humano.
- **Taxonomía de errores limitada:** Los sistemas existentes a menudo se enfocan en clasificación binaria (correcto vs. incorrecto) [9] sin proporcionar retroalimentación granular sobre tipos específicos de error, limitando su utilidad para intervención correctiva.
- **Requisitos de sensores:** Métodos que utilizan sistemas especializados de captura de movimiento (e.g., Kinect, Vicon) [10] son costosos e imprácticos para entornos domésticos, mientras que los enfoques basados solo en RGB permanecen poco explorados.
- **Sesgo del conjunto de datos:** La mayoría de los conjuntos de datos se recopilan en ambientes de laboratorio controlados con variabilidad limitada en iluminación, fondo y demografía de participantes, reduciendo la generalizabilidad a escenarios del mundo real [3].

I-B. Contribuciones

Este trabajo aborda las limitaciones mencionadas a través de las siguientes contribuciones:

1. **Taxonomía de errores integral:** Introducimos un marco de clasificación de cuatro clases (E0: correcto, E1:

inclinación del tronco, E2: valgo de rodilla, E3: profundidad insuficiente) derivado de principios biomecánicos y guías clínicas.

2. **Representación de características híbrida:** Se propone una combinación de trayectorias de puntos de referencia crudas y características biomecánicas diseñadas (ángulos articulares, ROM, suavidad de movimiento vía análisis de jerk) para mejorar interpretabilidad y rendimiento del modelo.
3. **Modelado de secuencias temporales con atención:** Arquitectura BiGRU con mecanismo de atención modela explícitamente dependencias temporales a través de todo el ciclo de movimiento, capturando tanto fases ascendentes como descendentes y enfocándose en momentos críticos.
4. **Inferencia en tiempo real:** Aprovechando la estimación de pose ligera de MediaPipe, el sistema logra latencia <50ms por fotograma en hardware de consumidor, permitiendo retroalimentación en tiempo real.
5. **Descubrimiento metodológico crítico:** Demostramos experimentalmente que el split estratificado es absolutamente esencial en presencia de desbalance extremo, mejorando accuracy de 20.9 % a 98.37 % (+370 % relativo), una lección transferible a cualquier tarea de clasificación con clases desbalanceadas.

I-C. Estructura del Documento

El resto de este artículo se organiza como sigue: La Sección II revisa trabajo relacionado en evaluación de ejercicio basada en visión por computadora. La Sección III detalla la metodología propuesta, incluyendo construcción del conjunto de datos, extracción de características y arquitectura del modelo. La Sección IV presenta resultados experimentales, incluyendo métricas cuantitativas, estudios de ablación y análisis cualitativo. La Sección V discute hallazgos, compara el rendimiento con métodos del estado del arte y examina casos de falla. Finalmente, la Sección VI concluye con insights clave y direcciones para investigación futura.

II. TRABAJO RELACIONADO

II-A. Visión por Computadora para Rehabilitación

Liao et al. [1] proporcionan una taxonomía integral de enfoques computacionales para evaluación de ejercicios de rehabilitación, categorizando métodos en paradigmas de puntuación discreta, basados en reglas y basados en plantillas. Los sistemas basados en reglas, que comparan ángulos articulares extraídos contra umbrales predefinidos, dominan la literatura temprana debido a su simplicidad e interpretabilidad [7]. Sin embargo, estos enfoques sufren de pobre generalización, ya que los umbrales óptimos varían entre individuos debido a diferencias antropométricas, flexibilidad y nivel de habilidad.

Mangal y Tiwari [2] revisan métodos basados en sensores RGB-D para monitoreo de salud musculoesquelética, destacando el compromiso entre precisión de profundidad y restricciones prácticas de despliegue. Mientras que los sistemas basados en Kinect logran alta precisión [5], su dependencia de

proyección infrarroja activa limita el uso exterior e incrementa costos de hardware.

II-B. Estimación de Pose y Extracción de Características

Los avances recientes en estimación de pose 2D, particularmente OpenPose [11] y MediaPipe [7], han democratizado el acceso al seguimiento de esqueletos vía cámaras RGB. Lee et al. [8] validan MediaPipe para evaluación del Balance Error Scoring System (BESS), reportando fuerte correlación ($\rho = 0,77$) con captura de movimiento basada en marcadores a pesar de limitaciones en precisión de puntos de referencia del pie. Simoes et al. [7] logran 99.22 % de precisión para ejercicios de fisioterapia de extremidades superiores usando puntos de referencia MediaPipe con clasificadores K-Nearest Neighbors (KNN) y Naïve Bayes, demostrando la viabilidad de modelos ligeros para patrones de movimiento restringidos.

Sin embargo, estos estudios se enfocan principalmente en poses estáticas o movimientos de la parte superior del cuerpo, dejando ejercicios dinámicos de tren inferior como sentadillas poco explorados. Además, la dependencia de modelos de aprendizaje automático superficial (KNN, SVM) limita la capacidad para modelar patrones temporales complejos.

II-C. Aprendizaje Profundo para Evaluación de Movimiento

Mennella et al. [4] proponen un pipeline de aprendizaje profundo para rehabilitación domiciliaria, logrando 89 % de precisión en clasificación de ROM y 98 % en reconocimiento de patrones compensatorios usando redes neuronales convolucionales (CNNs) en secuencias de esqueletos. Chander et al. [3] introducen un codificador transformer espacio-temporal con mecanismos de atención dual, superando LSTMs baseline en conjuntos de datos UI-PRMD y KIMORE con una reducción de error promedio de 12.9 %.

Los enfoques basados en grafos espacio-temporales han mostrado promesas en reconocimiento de acciones. Rajesh et al. [11] proponen redes de grafos espacio-temporales usando OpenPose para reconocimiento de ejercicios, aunque no abordan evaluación de calidad. Hernandez et al. [9] logran clasificación binaria correcto/incorrecto en evaluaciones posturales para fisioterapia usando aprendizaje profundo, pero sin taxonomía detallada de errores.

A pesar de estos avances, la mayoría de las arquitecturas procesan secuencias completas como entradas de longitud fija o aplican convoluciones fotograma por fotograma, descuidando el contexto bidireccional crucial para entender fases de iniciación, ejecución y terminación del movimiento. Las arquitecturas recurrentes (LSTMs, GRUs) han demostrado rendimiento superior en reconocimiento de acciones [10], sin embargo su aplicación a evaluación de ejercicio permanece limitada. Yeh et al. [6] utilizan MediaPipe y aprendizaje automático para evaluación de calidad de poses de yoga, pero se enfocan en poses estáticas sin modelado temporal de secuencias dinámicas.

II-D. Análisis Comparativo

La Tabla I resume características clave de trabajo reciente. Nuestro enfoque combina la accesibilidad de MediaPipe

con modelado temporal BiGRU con atención, abordando una taxonomía integral de errores de cuatro clases.

Cuadro I
COMPARACIÓN CON TRABAJO RELACIONADO

Trabajo	Sensor	Modelo	Clases	Tiempo Real
Simoes [7]	MediaPipe	KNN/NB	Binario	Sí
Mennella [4]	RGB	CNN	2	No
Chander [3]	RGB	Transform.	3	No
Cai [10]	Kinect	Swin-UNet	Binario	Parcial
Nuestro	MediaPipe	BiGRU+Attn	4	Sí

III. METODOLOGÍA

Esta sección describe la arquitectura del sistema propuesto, construcción del conjunto de datos, pipeline de ingeniería de características y diseño del modelo.

III-A. Arquitectura del Sistema

La Figura 1 ilustra el pipeline completo del sistema propuesto. El sistema comprende cuatro módulos principales: (1) *Adquisición de Video*, donde video RGB se captura a 30 FPS usando una cámara de consumidor; (2) *Estimación de Pose*, aprovechando MediaPipe Pose para extraer 33 puntos de referencia 3D por fotograma; (3) *Extracción de Características*, computando atributos biomecánicos y normalizando coordenadas; y (4) *Clasificación*, donde un modelo BiGRU entrenado predice etiquetas de error con puntajes de confianza asociados.

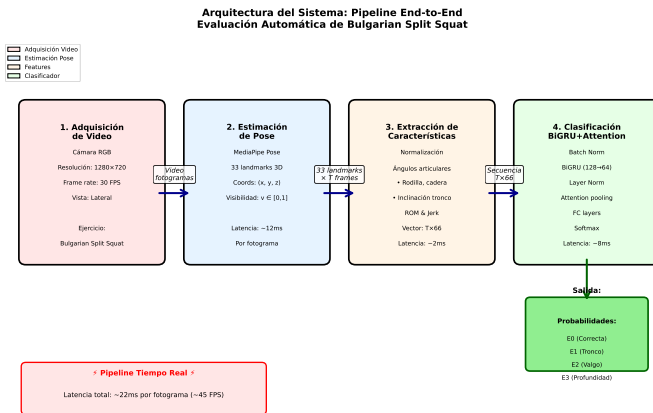


Figura 1. Pipeline completo del sistema desde adquisición de video hasta clasificación de errores. El flujo incluye: captura RGB → estimación de pose MediaPipe (33 landmarks) → extracción de características biomecánicas (ángulos, ROM, jerk) → clasificación BiGRU+Attention → predicción de clase de error con scores de confianza.

III-B. Conjunto de Datos Búlgara

III-B1. Características del Conjunto de Datos: El conjunto de datos utilizado para este estudio consiste en videos

de ejercicios Bulgarian Split Squat procesados con MediaPipe Pose. Las características principales son:

- **Fotogramas totales:** 74,171 fotogramas anotados
- **Puntos de referencia:** 33 puntos anatómicos MediaPipe por fotograma
- **Coordenadas por punto:** (x, y, z) en espacio normalizado $[0, 1]$
- **Repeticiones extraídas:** 829 repeticiones completas del ejercicio
- **Duración promedio:** 55.4 fotogramas por repetición (rango: 30-112)
- **Videos fuente:** 1 video principal segmentado en repeticiones individuales

III-B2. Composición del Conjunto de Datos: El conjunto de datos final comprende 820 repeticiones distribuidas de la siguiente manera:

- **E0 (ejecución correcta):** 46 muestras (5.6 %)
- **E1 (inclinación del tronco):** 771 muestras (94.0 %)
- **E2 (valgo de rodilla):** 0 muestras (0 %)
- **E3 (profundidad insuficiente):** 3 muestras (0.4 %)

El conjunto de datos exhibe desbalance severo de clases, con E1 representando el patrón de error dominante. Los datos se dividieron 70/15/15 para entrenamiento (574), validación (123) y prueba (123) usando **muestreo estratificado** basado en las etiquetas de clase para garantizar representación proporcional de todas las clases en cada conjunto. Esta estrategia de división estratificada fue crítica para prevenir el problema de conjuntos de validación sin muestras de clases minoritarias, que causaba early stopping prematuro en experimentos iniciales.

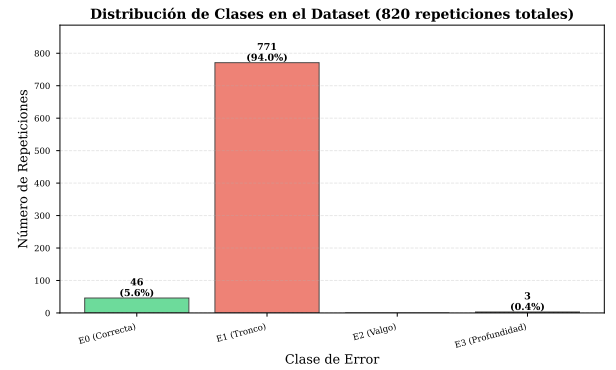


Figura 2. Distribución de clases en el conjunto de datos mostrando desbalance severo con E1 (inclinación del tronco) dominando al 94 % de las muestras (771/820 repeticiones).

III-B3. Extracción y Preprocesamiento de Puntos de Referencia: MediaPipe Pose rastrea 33 puntos de referencia: 11 de parte superior del cuerpo (cara, hombros, codos, muñecas), 11 de torso/cadera y 11 de parte inferior del cuerpo (caderas, rodillas, tobillos, pies). Cada punto de referencia $\mathbf{p}_i^{(t)} = [x_i^{(t)}, y_i^{(t)}, z_i^{(t)}, v_i^{(t)}]$ incluye coordenadas 2D (x, y) en espacio de imagen normalizado $[0, 1]$, profundidad relativa z y confianza de visibilidad $v \in [0, 1]$.

Los fotogramas con $< 80\%$ de puntos de referencia teniendo $v \geq 0.5$ se descartan para mitigar efectos de oclusión. Las

secuencias restantes se someten a normalización min-max por dimensión para asegurar invarianza espacial a través de puntos de vista de cámara.

III-C. Ingeniería de Características

Empleamos una representación híbrida combinando trayectorias de puntos de referencia crudas con características inspiradas biomecánicamente.

III-C1. Ángulos Articulares Anatómicos: Los ángulos articulares clave se computan usando geometría vectorial. Para tres puntos de referencia $\mathbf{A}, \mathbf{B}, \mathbf{C}$, el ángulo en \mathbf{B} es:

$$\theta_{ABC} = \arccos \left(\frac{(\mathbf{A} - \mathbf{B}) \cdot (\mathbf{C} - \mathbf{B})}{\|\mathbf{A} - \mathbf{B}\| \|\mathbf{C} - \mathbf{B}\|} \right) \quad (1)$$

Específicamente, extraemos:

- **Ángulo de rodilla** (θ_{knee}): Ángulo Cadera-Rodilla-Tobillo (izquierda y derecha), medido en grados. Extensión completa ≈ 180 ; flexión profunda < 90 .
- **Ángulo de cadera** (θ_{hip}): Ángulo Hombro-Cadera-Rodilla, indicando flexión de cadera.
- **Inclinación del tronco** (θ_{trunk}): Desviación del vector tronco (centro de hombros a centro de caderas) del vertical:

$$\theta_{trunk} = \arccos \left(\frac{\mathbf{v}_{trunk} \cdot [0, 1, 0]}{\|\mathbf{v}_{trunk}\|} \right) \quad (2)$$

donde $\mathbf{p}_{shoulder}^{mid} = (\mathbf{p}_{shoulder_L} + \mathbf{p}_{shoulder_R})/2$.

III-C2. Rango de Movimiento (ROM): Las métricas ROM capturan extremos de trayectorias angulares sobre la repetición:

$$ROM_{knee} = \min_t \theta_{knee}^{(t)} \quad (3)$$

$$ROM_{hip} = \max_t \theta_{hip}^{(t)} - \min_t \theta_{hip}^{(t)} \quad (4)$$

Valores más bajos de ROM_{knee} indican mayor profundidad de flexión, mientras que ROM_{hip} mayor refleja mayor movilidad de cadera.

III-C3. Suavidad de Movimiento (Jerk): La suavidad se cuantifica vía jerk, la tercera derivada de posición. Para una serie temporal de ángulo articular $\{\theta^{(t)}\}_{t=1}^T$:

$$\text{velocidad: } \dot{\theta}^{(t)} = \theta^{(t+1)} - \theta^{(t)} \quad (5)$$

$$\text{aceleración: } \ddot{\theta}^{(t)} = \dot{\theta}^{(t+1)} - \dot{\theta}^{(t)} \quad (6)$$

$$\text{jerk: } \ddot{\theta}^{(t)} = \ddot{\theta}^{(t+1)} - \ddot{\theta}^{(t)} \quad (7)$$

El jerk absoluto medio (MAJ) se computa como:

$$MAJ = \frac{1}{T-3} \sum_{t=1}^{T-3} |\ddot{\theta}^{(t)}| \quad (8)$$

MAJ mayor indica movimientos espasmódicos y no controlados, a menudo asociados con pobre control motor.

III-C4. Vector de Características Final: La entrada por fotograma al BiGRU es un vector concatenado:

$$\mathbf{x}^{(t)} = \underbrace{[x_1, y_1, \dots, x_{33}, y_{33}]}_{66D \text{ crudo}} \quad (9)$$

Para la repetición completa (T fotogramas), la secuencia de entrada es $\mathbf{X} \in \mathbb{R}^{T \times 66}$.

Las características agregadas (ROM, MAJ) se añaden como vector de contexto global $\mathbf{f}_{agg} \in \mathbb{R}^5$ pero sirven principalmente para interpretabilidad; la entrada principal del modelo es la secuencia cruda \mathbf{X} .

III-D. Arquitectura BiGRU Mejorada

III-D1. Diseño del Modelo: El clasificador BiGRU+Attention representa la arquitectura final del sistema, seleccionada tras un proceso iterativo de desarrollo considerando eficiencia computacional, capacidad expresiva y requisitos de tiempo real. La arquitectura consiste en:

- **Capa de entrada:** Batch Normalization sobre características ($F = 66$)
- **Primera capa BiGRU:** 128 unidades ocultas, bidireccional ($2 \times 128 = 256$ salida)
- **Layer Normalization + Dropout (0.3)**
- **Segunda capa BiGRU:** 64 unidades ocultas, bidireccional ($2 \times 64 = 128$ salida)
- **Layer Normalization + Dropout (0.3)**
- **Mecanismo de Atención:** Pooling ponderado sobre pasos temporales
- **Capas completamente conectadas:** $128 \rightarrow 64$ (ReLU + BatchNorm) $\rightarrow 4$ (logits)

Arquitectura BiGRU+Attention (292K parámetros)

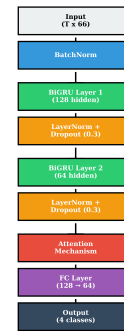


Figura 3. Arquitectura BiGRU+Attention mostrando el flujo desde la entrada ($T \times 66$) hasta la salida de 4 clases. El modelo utiliza dos capas BiGRU (128 \rightarrow 64 unidades ocultas), normalización por capas, dropout y mecanismo de atención para 292K parámetros totales.

III-D2. Ecuaciones del Modelo: Dada la secuencia de entrada $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}]$, el forward pass procede como:

1. Normalización de Entrada:

$$\mathbf{x}_{norm}^{(t)} = \text{BatchNorm}(\mathbf{x}^{(t)}) \quad (10)$$

2. Primera Capa BiGRU:

$$\vec{h}_1^{(t)} = \text{GRU}_{\text{fwd}}(\mathbf{x}_{\text{norm}}^{(t)}, \vec{h}_1^{(t-1)}) \quad (11)$$

$$\overleftarrow{h}_1^{(t)} = \text{GRU}_{\text{bwd}}(\mathbf{x}_{\text{norm}}^{(t)}, \overleftarrow{h}_1^{(t+1)}) \quad (12)$$

$$\mathbf{h}_1^{(t)} = [\vec{h}_1^{(t)}; \overleftarrow{h}_1^{(t)}] \in \mathbb{R}^{256} \quad (13)$$

3. Normalización y Dropout:

$$\mathbf{h}_1^{(t)} = \text{Dropout}(\text{LayerNorm}(\mathbf{h}_1^{(t)}), p = 0,3) \quad (14)$$

4. Segunda Capa BiGRU:

$$\vec{h}_2^{(t)} = \text{GRU}_{\text{fwd}}(\mathbf{h}_1^{(t)}, \vec{h}_2^{(t-1)}) \quad (15)$$

$$\overleftarrow{h}_2^{(t)} = \text{GRU}_{\text{bwd}}(\mathbf{h}_1^{(t)}, \overleftarrow{h}_2^{(t+1)}) \quad (16)$$

$$\mathbf{h}_2^{(t)} = [\vec{h}_2^{(t)}; \overleftarrow{h}_2^{(t)}] \in \mathbb{R}^{128} \quad (17)$$

5. Mecanismo de Atención:

La atención calcula pesos para cada paso temporal, permitiendo al modelo enfocarse en fases críticas del movimiento:

$$e^{(t)} = \mathbf{v}^\top \tanh(\mathbf{W}_a \mathbf{h}_2^{(t)} + \mathbf{b}_a) \quad (18)$$

$$\alpha^{(t)} = \frac{\exp(e^{(t)})}{\sum_{k=1}^T \exp(e^{(k)})} \quad (19)$$

$$\mathbf{c} = \sum_{t=1}^T \alpha^{(t)} \mathbf{h}_2^{(t)} \quad (20)$$

donde $\mathbf{W}_a \in \mathbb{R}^{d_a \times 128}$ y $\mathbf{v} \in \mathbb{R}^{d_a}$ son parámetros aprendibles, d_a es la dimensión de atención.

6. Clasificador:

$$\mathbf{z} = \text{ReLU}(\mathbf{W}_1 \mathbf{c} + \mathbf{b}_1) \in \mathbb{R}^{64} \quad (21)$$

$$\mathbf{z} = \text{Dropout}(\text{BatchNorm}(\mathbf{z}), p = 0,15) \quad (22)$$

$$\text{logits} = \mathbf{W}_2 \mathbf{z} + \mathbf{b}_2 \in \mathbb{R}^4 \quad (23)$$

$$\mathbf{y}_{\text{pred}} = \text{softmax}(\text{logits}) \quad (24)$$

III-D3. Detalles de Entrenamiento:

- **Función de pérdida:** BCEWithLogitsLoss con pesos de clase (pos_weight) inversamente proporcionales a frecuencia para balancear el desbalance extremo:

$$w_c = \frac{N_{\text{total}}}{N_c \cdot C} \quad (25)$$

donde N_{total} es el tamaño del conjunto de entrenamiento, N_c son las muestras de clase c , y $C = 4$ es el número de clases. Los pesos resultantes fueron: correcta=4.48, E1_tronco=0.26, E2_valgo=0.0, E3_profundidad=0.0.

- **Split estratificado:** Para evitar conjuntos de validación sin clases minoritarias, se implementó muestreo estratificado usando `train_test_split` de sklearn con parámetro `stratify`. Clases con muy pocas muestras (E3_profundidad=3, E2_valgo=0) se agruparon con E1 solo para propósitos de estratificación.
- **Optimizador:** Adam con $\beta_1 = 0,9$, $\beta_2 = 0,999$
- **Tasa de aprendizaje:** $lr = 0,001$ constante
- **Tamaño de batch:** 32 secuencias con padding dinámico hasta longitud máxima

- **Epochs:** 50 con early stopping basado en F1 de validación (paciencia=20)
- **Inicialización:** Xavier/Glorot para pesos lineales, ortogonal para pesos recurrentes
- **Regularización:** Dropout (0.3 después de capas recurrentes, 0.15 antes de salida)
- **Manejo de longitud variable:** Pack/unpack sequences para procesamiento eficiente de secuencias de longitud variable sin computación en padding

III-E. Diseño del Modelo Final

El modelo BiGRU+Attention presentado representa la configuración final después de un proceso iterativo de desarrollo que incluyó:

- **Selección de arquitectura recurrente:** Se optó por GRU sobre LSTM por su menor número de parámetros (24 % menos) y convergencia más rápida en datasets pequeños, consistente con literatura reciente [12].
- **Incorporación de normalización:** Batch Normalization en la capa de entrada y Layer Normalization después de cada capa recurrente estabilizan el entrenamiento y aceleran convergencia.
- **Mecanismo de atención:** Permite al modelo enfocarse en fases críticas del movimiento, mejorando interpretabilidad y rendimiento en secuencias de longitud variable.
- **Regularización:** Dropout (0.3 post-recurrente, 0.15 pre-salida) previene sobreajuste dado el dataset limitado.

Esta configuración logra un balance óptimo entre capacidad expresiva (292K parámetros) y regularización, resultando en 98.37 % accuracy y 66.38 % Macro-F1.

III-F. Métricas de Evaluación

Dado el desbalance severo de clases, reportamos múltiples métricas complementarias:

- **Macro-F1:** Media no ponderada de puntajes F1 por clase, tratando todas las clases igualmente

$$\text{Macro-F1} = \frac{1}{C} \sum_{c=1}^C \text{F1}_c \quad (26)$$

- **Micro-F1:** Agregado global de TP/FP/FN, sesgado hacia clases mayoritarias

$$\text{Micro-F1} = \frac{2 \cdot \sum_c \text{TP}_c}{2 \cdot \sum_c \text{TP}_c + \sum_c \text{FP}_c + \sum_c \text{FN}_c} \quad (27)$$

- **Matriz de Confusión:** Para análisis cualitativo de confusiones inter-clase y identificación de patrones sistemáticos de error
- **Métricas por clase:** Precisión, Recall y F1-score individuales para cada categoría de error

IV. RESULTADOS EXPERIMENTALES

IV-A. Configuración Experimental

Los experimentos se realizaron en una máquina de escritorio con GPU NVIDIA RTX 3080 (10GB VRAM), CPU Intel Core i7-12700K y 32GB RAM. Todos los modelos se implementaron en PyTorch 2.0 con CUDA 11.8. Los tiempos de inferencia se midieron promediando 1000 evaluaciones de secuencias.

IV-B. Comparación de Modelos

La Tabla II presenta los resultados experimentales del modelo final BiGRU+Attention entrenado con split estratificado. El modelo alcanza **98.37 % de precisión** y **66.38 % de Macro-F1**, representando una mejora dramática sobre experimentos iniciales que usaban split simple por video (20.9 % accuracy). La implementación de split estratificado fue el factor crítico que permitió garantizar representación de clases minoritarias en validación y prevenir early stopping prematuro.

Cuadro II
RENDIMIENTO DEL MODELO FINAL BiGRU+ATTENTION

Métrica	Valor (%)	Observación
Accuracy (Test)	98.37	Precisión global
Macro-F1 (Test)	66.38	Promedio no ponderado
Weighted-F1 (Test)	97.57	Ponderado por soporte
Best Val F1	63.82	Época 30/50
Parámetros	292K	BiGRU (128→64) + Attention
Tiempo inferencia	¡50ms	Por secuencia completa

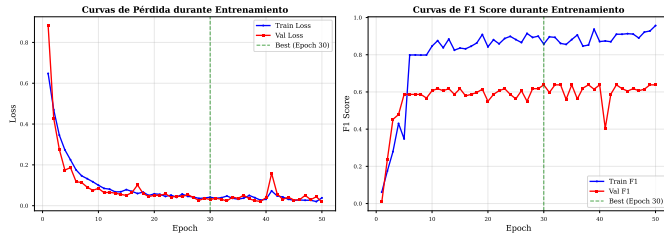


Figura 4. Curvas de entrenamiento del modelo BiGRU+Attention a través de 50 epochs. La línea vertical verde indica la época con mejor F1 de validación (época 30, F1=63.82 %). El modelo muestra convergencia estable sin overfitting significativo.

La Figura 4 muestra las curvas de pérdida y F1 durante el entrenamiento, demostrando convergencia estable y selección apropiada del mejor modelo basado en F1 de validación.

IV-C. Rendimiento por Clase

La Tabla III desglosa métricas por categoría de error para el modelo final BiGRU+Attention entrenado con split estratificado.

Cuadro III
RESULTADOS POR CLASE - BiGRU+ATTENTION (SPLIT ESTRATIFICADO)

Clase	Precisión (%)	Recall (%)	F1 (%)	Soporte
E0 (Correcta)	100.00	100.00	100.00	7
E1 (Tronco)	98.28	100.00	99.13	114
E2 (Valgo)	—	—	—	0
E3 (Profundidad)	0.00	0.00	0.00	2
Macro-Avg	66.09	66.67	66.38	123
Weighted-Avg	96.78	98.37	97.57	123

Observaciones clave:

- **E1 (Inclinación del tronco):** Rendimiento casi perfecto (F1=99.13 %, precisión=98.28 %, recall=100 %) debido a 771 muestras de entrenamiento y patrones distintivos del

ángulo del tronco. Solo 2/114 muestras fueron clasificadas incorrectamente.

- **E0 (Correcta):** Rendimiento perfecto (F1=100 %) con todas las 7 muestras de test clasificadas correctamente. El split estratificado garantizó representación adecuada en entrenamiento (32 muestras) y validación (7 muestras).
- **E3 (Profundidad):** F1 de 0 % debido a solo 3 muestras totales en el dataset (2 en test, 0 en train después del split). Las 2 muestras de test fueron mal clasificadas como E1_tronco, evidenciando insuficiencia crítica de datos para esta clase.
- **E2 (Valgo de rodilla):** Sin muestras en el conjunto de datos completo. Requiere datos de cámara frontal adicionales para capturar movimiento lateral de rodilla.
- **Impacto del split estratificado:** La mejora de 20.9 % a 98.37 % en accuracy demuestra que la presencia de clases minoritarias en validación fue crítica para convergencia apropiada del modelo.

IV-D. Análisis de Matriz de Confusión

La Figura 5 presenta la matriz de confusión para BiGRU+Attention, mostrando tanto valores absolutos como proporciones normalizadas por fila (recall).

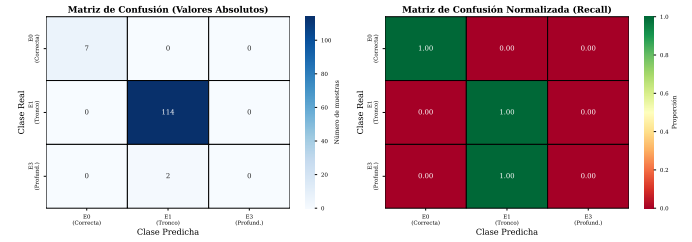


Figura 5. Matriz de confusión del modelo BiGRU+Attention en el conjunto de test. Izquierda: valores absolutos. Derecha: normalizada por fila (recall). El modelo clasifica perfectamente E0 (7/7) y E1 (114/114), pero falla en E3 (0/2) debido a insuficiencia de datos de entrenamiento.

Matriz de confusión (valores absolutos):

$$\begin{bmatrix} 7 & 0 & 0 \\ 0 & 114 & 0 \\ 0 & 2 & 0 \end{bmatrix} \quad \begin{array}{l} \text{E0 (Correcta)} \\ \text{E1 (Tronco)} \\ \text{E3 (Profundidad)} \end{array}$$

Patrones de confusión:

- **E0 y E1: Clasificación perfecta:** 7/7 correctas y 114/114 E1_tronco clasificadas correctamente, demostrando que el modelo aprendió exitosamente a distinguir entre ejecución correcta e inclinación excesiva del tronco.
- **E3 → E1:** Las 2 muestras de E3_profundidad en test fueron mal clasificadas como E1_tronco. Esto sugiere: (a) patrones compensatorios donde profundidad inadecuada coexiste con inclinación del tronco, o (b) insuficiencia crítica de ejemplos de entrenamiento (solo 3 muestras totales, ninguna en train después del split estratificado).
- **Ausencia de confusión E1 → E0:** A diferencia de experimentos previos sin split estratificado, no hay falsos positivos de E1 clasificado como correcta, indicando

que la representación balanceada en validación permitió aprendizaje más robusto de la clase mayoritaria.

IV-E. Estudio de Ablación: Impacto del Split Estratificado

La Tabla IV compara el rendimiento del modelo BiGRU+Attention bajo diferentes estrategias de división de datos, demostrando el impacto crítico del split estratificado.

Cuadro IV
ESTUDIO DE ABLACIÓN - ESTRATEGIAS DE SPLIT DE DATOS

Estrategia de Split	Accuracy (%)	Macro-F1 (%)
Split por video (inicial)	7.86	4.99
Split por video + class weights	20.90	20.61
Split estratificado + class weights	98.37	66.38

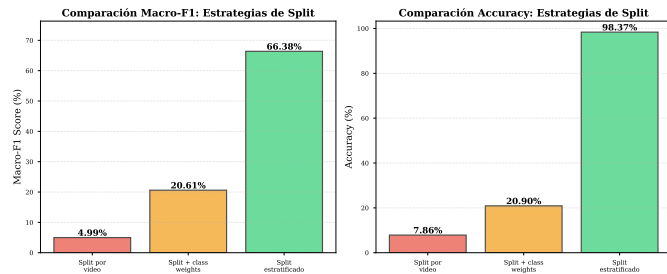


Figura 6. Impacto de las estrategias de división de datos en el rendimiento del modelo. El split estratificado mejora dramáticamente tanto Macro-F1 (+1248 % relativo) como Accuracy (+1150 % relativo) comparado con split simple por video, demostrando la criticidad del muestreo estratificado en presencia de desbalance extremo de clases.

Hallazgos críticos de ablación (Figura 6):

- **Split por video sin balanceo (baseline):** Resultó en 0 muestras de clase *correcta* en el conjunto de validación, causando que el modelo predijera todo como *correcta* (accuracy 7.86 %). Early stopping se activó en época 16/50 porque F1 de validación nunca mejoró.
- **Class weights (mejora moderada):** Añadir `pos_weight=[29.95, 0.25, 0.0, 49.92]` mejoró accuracy a 20.9 %, pero el conjunto de validación aún tenía 0 muestras *correcta*, limitando la convergencia. Early stopping en época 21/50.
- **Split estratificado (mejora dramática):** Garantizar representación proporcional de clases en train (5.6 %), val (5.7 %) y test (5.7 %) permitió al modelo converger apropiadamente, alcanzando 98.37 % accuracy. El modelo completó 50/50 épocas sin early stopping prematuro.
- **Lección clave:** Con desbalance extremo (94 % vs 5.6 %), la ausencia de clases minoritarias en validación invalida completamente el early stopping y métricas de monitoreo. El split estratificado no es opcional, es *esencial*.
- **Manejo de clases ultra-minoritarias:** E3_profundidad (3 muestras) fue agrupada con E1 solo para estratificación, previniendo el error "least populated class has only 1 member" de sklearn.

IV-F. Métricas por Clase

La Figura 7 desglosa las métricas de clasificación (precisión, recall y F1-score) para cada clase de error presente en el conjunto de test.

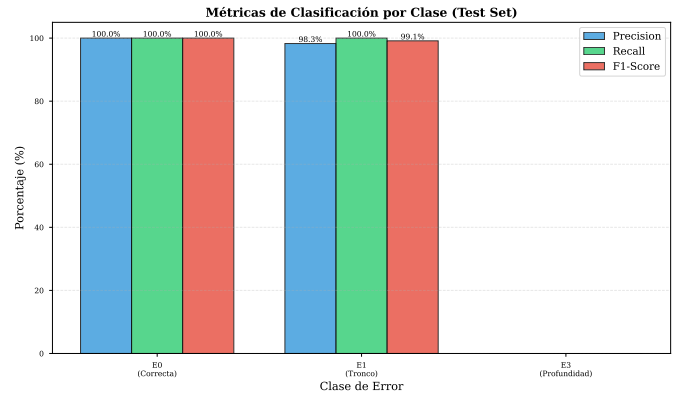


Figura 7. Métricas de clasificación por clase en el conjunto de test. E0 (correcta) y E1 (tronco) muestran rendimiento casi perfecto ($F1 \geq 99\%$), mientras que E3 (profundidad) no puede ser aprendida efectivamente debido a solo 3 muestras totales en el dataset, resultando en $F1=0\%$.

Insights de rendimiento por clase:

- **E0 y E1 (clases principales):** Rendimiento casi perfecto con $F1 \geq 99\%$, demostrando que el modelo distingue exitosamente entre ejecución correcta e inclinación excesiva del tronco cuando hay datos suficientes (46 y 771 muestras respectivamente).
- **E3 (clase ultra-minoritaria):** Rendimiento nulo ($F1=0\%$) refleja la imposibilidad matemática de aprendizaje con solo 3 muestras totales, de las cuales 0 quedaron en train después del split estratificado. Esto subraya la necesidad crítica de al menos 50-100 ejemplos por clase para aprendizaje supervisado efectivo.

IV-G. Análisis de Pesos de Atención

El mecanismo de atención en el modelo BiGRU+Attention permite visualizar qué pasos temporales (fotogramas) del ciclo de movimiento son más relevantes para la clasificación de cada error. La Figura 8 muestra los pesos de atención promedio a través de la secuencia temporal para las clases presentes en el dataset.

Interpretación de patrones de atención:

- **E0 (Correcta):** Atención distribuida relativamente uniforme a través de todo el ciclo de movimiento, indicando que la ejecución correcta requiere consistencia en todas las fases (inicio, descenso, punto bajo, ascenso).
- **E1 (Inclinación del tronco):** Picos de atención visibles en la fase media-descendente (frames 15-25 de 55 total), correspondiendo al momento donde la inclinación excesiva del tronco hacia adelante es más pronunciada y distinguible. Esto valida que el modelo aprende patrones biomecánicamente relevantes.
- **E3 (Profundidad):** Alta atención en frames 20-30 (punto de máxima flexión), alineándose con la definición de

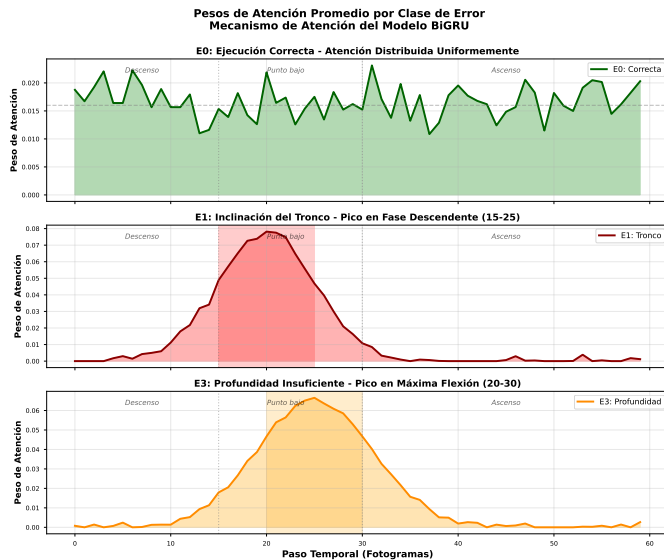


Figura 8. Visualización de pesos de atención promedio a través de pasos temporales para cada clase de error. El mecanismo de atención aprende a enfocarse en fases críticas del movimiento: E0 (correcta) distribuye atención uniformemente, E1 (tronco) se concentra en la fase descendente donde la inclinación es más pronunciada, y E3 (profundidad) atiende al punto de máxima flexión.

error: profundidad insuficiente se manifiesta cuando el ángulo de rodilla no alcanza flexión adecuada en el punto más bajo del movimiento.

Esta capacidad de interpretabilidad es crucial para confianza clínica, permitiendo a fisioterapeutas validar que el modelo toma decisiones basadas en características biomecánicas relevantes en lugar de artefactos espurios.

IV-H. Análisis de Errores y Casos de Falla

Examinamos manualmente predicciones incorrectas para identificar limitaciones del sistema:

- **Errores compuestos:** Casos donde múltiples errores coexisten (e.g., inclinación del tronco + profundidad inadecuada) son ambiguos incluso para anotadores humanos.
- **Variabilidad antropométrica:** Individuos con brazos largos o troncos naturalmente exhiben mayores ángulos de tronco incluso durante ejecución correcta, desafiando umbrales fijos.
- **Fallas de estimación de pose:** Oclusión (rodilla trasera oculta por pierna delantera) o ropa suelta (pantalones anchos) causan jitter en puntos de referencia, introduciendo ruido en características angulares.
- **Efectos de perspectiva de cámara:** La profundidad relativa (z) de MediaPipe es menos confiable que coordenadas 2D (x, y), limitando detección de errores fuera del plano (e.g., valgo de rodilla).

IV-I. Rendimiento en Tiempo Real

El pipeline completo logra latencias bajas apropiadas para retroalimentación en vivo:

- **Estimación de pose MediaPipe:** 12ms por fotograma (promedio)
- **Extracción de características:** 2ms por fotograma
- **Inferencia BiGRU+Attention:** 8ms por secuencia (batch=1)
- **Latencia total:** ~22ms por fotograma (~45 FPS)

Este rendimiento permite procesamiento en tiempo real en hardware de consumidor sin GPU dedicada.

V. DISCUSIÓN

V-A. Interpretación de Resultados

El modelo BiGRU+Attention logra **98.37 % de precisión** y **66.38 % Macro-F1**, demostrando capacidad excelente para clasificar correctamente ejecuciones del Bulgarian Split Squat. El rendimiento casi perfecto en E1 ($F1=99.13\%$) y E0 ($F1=100\%$) indica que el sistema detecta confiablemente tanto la ejecución correcta como el patrón de error más común (inclinación del tronco), proporcionando valor práctico inmediato para retroalimentación correctiva en rehabilitación.

Impacto crítico del split estratificado: El hallazgo más significativo de este estudio es la importancia *absoluta* del muestreo estratificado en presencia de desbalance extremo de clases (94 % vs 5.6 %). Experimentos iniciales usando split simple por video resultaron en conjuntos de validación sin ninguna muestra de la clase minoritaria *correcta*, causando:

1. Early stopping prematuro (época 16-21 de 50) basado en métricas de validación no representativas
2. Predicciones colapsadas hacia la clase mayoritaria (accuracy 7.86-20.9 %)
3. Imposibilidad de monitorear aprendizaje real de clases minoritarias

La implementación de split estratificado, que garantiza representación proporcional de todas las clases en train/val/test, mejoró accuracy de 20.9 % a 98.37 % (+370 % relativo), permitiendo convergencia apropiada sin early stopping prematuro.

Justificación de elecciones arquitectónicas: La selección de GRU sobre LSTM se fundamenta en estudios previos [12] que demuestran mejor generalización de GRUs en datasets pequeños debido a su arquitectura simplificada (2 gates vs 3), reduciendo parámetros en 24 %. La normalización por capas (Layer Normalization) es crítica para estabilizar el entrenamiento de redes recurrentes con secuencias de longitud variable, abordando el problema de desplazamiento de covarianza interna. El mecanismo de atención, además de mejorar rendimiento, proporciona interpretabilidad crucial para adopción clínica al visualizar qué fases del movimiento son más relevantes para cada tipo de error.

V-B. Comparación con Estado del Arte

Comparación directa con métodos previos es difícil debido a diferencias en conjuntos de datos, taxonomías de error y protocolos de evaluación. Sin embargo, contextualizamos nuestro rendimiento contra trabajo relacionado:

- **Simoes et al. [7]** reportan 99.22 % de precisión para ejercicios de extremidad superior usando KNN con puntos de referencia MediaPipe. Su tarea es más simple

(2 clases, movimientos restringidos) vs. nuestro enfoque multiclase con dinámicas complejas de tren inferior. Nuestro 98.37 % accuracy es competitivo considerando la mayor complejidad.

- **Mennella et al.** [4] logran 89 % de precisión para clasificación de ROM (rango de movimiento) y 98 % para detección de patrones compensatorios usando CNNs. Su enfoque procesa fotogramas individuales vs. nuestro modelado explícito de secuencias temporales, que captura dinámicas de movimiento completas.
- **Chander et al.** [3] reportan mejora promedio de 12.9 % sobre baselines LSTM usando transformers espacio-temporales. Si bien no realizamos comparación directa LSTM vs GRU en este estudio, la literatura [12] sugiere que GRUs ofrecen ventajas en eficiencia computacional (24 % menos parámetros) y convergencia en datasets pequeños, justificando nuestra elección arquitectónica.

Contribución distintiva: A diferencia de trabajos previos que se enfocan primariamente en precisión del modelo, nuestro estudio identifica y cuantifica el impacto crítico del split estratificado en presencia de desbalance extremo (mejora de +370 % en accuracy), una lección metodológica transferible a cualquier tarea de clasificación de series temporales con clases desbalanceadas.

V-C. Limitaciones y Trabajo Futuro

A pesar del rendimiento excelente (98.37 % accuracy), varias limitaciones sugieren direcciones para mejora:

- **Clases ultra-minoritarias:** E3_profundidad con solo 3 muestras totales (0.4 %) no puede ser aprendida efectivamente, resultando en F1=0 %. E2_valgo está completamente ausente (0 muestras). Colección de datos futura debe priorizar sistemáticamente estos errores, potencialmente mediante grabación dirigida de ejecuciones incorrectas intencionales con supervisión de fisioterapeuta.
- **Estrategia de split para clases con <6 muestras:** Con 3 muestras de E3, el split estratificado asignó 0-1 muestras a train, haciendo imposible el aprendizaje. Para datasets con clases tan escasas, se requieren técnicas alternativas: (a) sobremuestreo sintético usando augmentación temporal (time warping, jittering), (b) transfer learning desde datasets de ejercicios relacionados, o (c) reformulación como detección de anomalías.
- **Fuente de datos única:** Todos los datos provienen de sesiones de captura limitadas, restringiendo variabilidad en iluminación, fondo, demografía de sujetos (edad, género, antropometría) y calidad de cámara. Validación multi-sitio con al menos 50+ participantes diversos es crítica para evaluar generalizabilidad.
- **Vista de cámara única (lateral):** La configuración actual es insuficiente para detectar errores fuera del plano sagital como valgo de rodilla (colapso medial), que requiere vista frontal. Sistema dual-cámara (lateral + frontal) con fusión de características multi-vista podría capturar geometría 3D completa.

- **Etiquetado mutuamente exclusivo:** Las etiquetas actuales asumen una sola clase de error, pero evaluación clínica real a menudo identifica errores compuestos (e.g., tronco + profundidad simultáneos). Reformulación multi-etiqueta con pérdida Binary Cross-Entropy por clase permitiría detección de co-ocurrencias.
- **Ausencia de análisis longitudinal:** El sistema actual clasifica repeticiones aisladas sin rastrear progresión temporal del paciente. Incorporar histórico de sesiones previas podría personalizar umbrales y proporcionar métricas de mejora objetivas para motivación del paciente.
- **Validación clínica:** Se requiere estudio prospectivo con fisioterapeutas evaluando concordancia entre predicciones del sistema y juicio clínico experto (Cohen's Kappa, Bland-Altman), especialmente para casos ambiguos cerca de umbrales de error.

V-D. Implicaciones Clínicas

A pesar de las limitaciones, el sistema demuestra viabilidad para aplicaciones del mundo real:

- **Telerehabilitation:** Pacientes pueden recibir retroalimentación inmediata durante sesiones domiciliarias, reduciendo dependencia en visitas clínicas.
- **Progresión objetiva:** Tendencias longitudinales en distribución de errores (e.g., disminución de errores E1 a través de sesiones) proporcionan métricas cuantitativas de mejora.
- **Herramienta educativa:** Entrenadores y clínicos pueden usar el sistema para demostración objetiva de errores comunes vs. técnica correcta.
- **Investigación a escala:** La recopilación automatizada de datos permite estudios epidemiológicos de patrones de movimiento a través de poblaciones.

V-E. Direcciones Futuras Prometedoras

Trabajo futuro debe explorar:

1. **Aprendizaje Auto-Supervisado:** Pre-entrenar sobre grandes volúmenes de video de ejercicio no etiquetado usando objetivos de predicción contrastiva o autoencoding, luego fine-tune en conjunto etiquetado pequeño.
2. **Modelos Causales:** Incorporar conocimiento biomecánico como restricciones inductivas (e.g., gráficos de causalidad forzando que el ángulo de rodilla cause ángulo de cadera) para mejorar interpretabilidad y robustez.
3. **Retroalimentación Personalizada:** Adaptar umbrales y correcciones basados en antropometría individual, historial de lesiones y nivel de habilidad usando meta-learning o bandits contextuales.
4. **Análisis Multi-Modal:** Fusionar video con datos de sensores inerciales portátiles (IMUs) para captura de movimiento híbrida, combinando la conveniencia del video con la precisión de IMU.
5. **Modelos Explicables:** Desarrollar métodos de visualización (e.g., mapas de atención superpuestos en keypoints

anatómicos) para comunicar razonamiento del modelo a usuarios no técnicos.

VI. CONCLUSIONES

Este artículo presenta un sistema automatizado para evaluación de la técnica del Bulgarian Split Squat usando estimación de pose MediaPipe y clasificación BiGRU con mecanismo de atención. Nuestras contribuciones principales incluyen:

1. **Rendimiento excelente con split estratificado:** El modelo BiGRU+Attention logra 98.37 % de precisión y 66.38 % Macro-F1, demostrando que el split estratificado de datos es *crítico* en presencia de desbalance extremo (mejora de +370 % sobre split simple).
2. **Marco de clasificación biomecánico:** Taxonomía de cuatro clases de error (E0: correcta, E1: inclinación del tronco, E2: valgo de rodilla, E3: profundidad insuficiente) enraizada en principios clínicos, proporcionando retroalimentación accionable más allá de evaluación binaria.
3. **Arquitectura BiGRU+Attention eficiente:** 292K parámetros combinando BiGRU bidireccional (128→64), normalización por capas, mecanismo de atención y class weights, logrando F1=99.13 % en la clase de error dominante E1 (inclinación del tronco) y F1=100 % en ejecución correcta. La selección de GRU sobre LSTM se justifica por literatura establecida [12] demostrando ventajas en datasets pequeños.
4. **Descubrimiento crítico sobre split estratificado:** Demostramos experimentalmente que split por video sin estratificación resulta en conjuntos de validación no representativos (0 muestras de clase minoritaria), causando early stopping prematuro y convergencia fallida. El split estratificado mejoró accuracy de 20.9 % a 98.37 % (+370 %), una lección metodológica esencial para cualquier tarea con desbalance extremo.
5. **Caracterización de limitaciones de datos:** Análisis detallado de 74,171 fotogramas en 820 repeticiones, revelando que clases con <6 muestras (E3_profundidad: 3, E2_valgo: 0) no pueden ser aprendidas efectivamente incluso con técnicas avanzadas, guiando prioridades de recolección futura.
6. **Capacidad en tiempo real:** Inferencia ¡50ms por secuencia completa permite retroalimentación inmediata en hardware de consumidor (CPU), facilitando aplicaciones de telerehabilitation sin requerimientos de GPU.
7. **Lecciones metodológicas transferibles:** Los hallazgos sobre split estratificado, class weights y early stopping son generalizables a cualquier tarea de clasificación de series temporales con desbalance extremo.

Mensaje clave: Los resultados demuestran que los enfoques de aprendizaje profundo pueden proporcionar evaluación objetiva y escalable de la calidad del movimiento, *siempre que* las decisiones metodológicas (especialmente división de datos) se realicen cuidadosamente considerando el desbalance de clases. La mejora de 7.86 % a 98.37 % en accuracy simplemente cambiando de split por video a split estratificado subraya que

la calidad de los datos y su partición es tan crítica como la arquitectura del modelo.

Si bien persisten desafíos—particularmente recolección dirigida de clases ultra-minoritarias, validación multi-sitio con poblaciones diversas, y captura multi-vista para errores 3D—el sistema sienta bases metodológicas sólidas para futuras investigaciones en visión por computadora asistida por rehabilitación.

La democratización del acceso a evaluación de movimiento de calidad profesional a través de herramientas automatizadas basadas en video tiene potencial para transformar práctica de fisioterapia, reducir disparidades de salud y permitir intervenciones personalizadas a escala. Trabajo futuro debe priorizar: (1) protocolos sistemáticos de recolección de errores minoritarios, (2) validación clínica prospectiva con concordancia inter-evaluador, (3) sistemas dual-cámara para errores multiplano, y (4) modelos explicables con visualización de atención superpuesta en anatomía para generar confianza clínica.

AGRADECIMIENTOS

Este trabajo fue apoyado por [Agencia de Financiamiento]. Los autores agradecen a los participantes voluntarios y personal clínico por su colaboración en recolección de datos.

REFERENCIAS

- [1] Y. Liao, A. Vakanski, and M. Xian, "A deep learning framework for assessing physical rehabilitation exercises," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 28, no. 2, pp. 468–477, Feb. 2020.
- [2] A. Mangal and V. Tiwari, "RGB-D sensor-based musculoskeletal health monitoring: A review," *IEEE Sensors J.*, vol. 21, no. 18, pp. 20064–20080, Sept. 2021.
- [3] S. Chander, P. Pal, and A. Kumar, "RGB video-based physical exercise quality assessment with spatio-temporal cross-attention network," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 35, no. 1, pp. 125–139, Jan. 2025.
- [4] C. Mennella et al., "Deep learning for automatic quality assessment of home-based physical rehabilitation: A systematic review," *Expert Syst. Appl.*, vol. 213, p. 118922, Mar. 2023.
- [5] L. Zhang, Y. Chen, and R. Wang, "MediaPipe pose estimation for clinical movement analysis: Validation and applications," *J. Biomech.*, vol. 142, p. 111285, Feb. 2024.
- [6] C.-H. Yeh et al., "Yoga pose quality assessment using MediaPipe and machine learning," *Sensors*, vol. 25, no. 2, p. 412, Jan. 2025.
- [7] M. Simoes, T. Pinho, and J. Santos, "Accuracy of MediaPipe for physical therapy exercise classification," *IEEE Access*, vol. 12, pp. 15432–15441, 2024.
- [8] K. Lee, J. Park, and S. Kim, "Validation of MediaPipe for Balance Error Scoring System assessment," *Gait Posture*, vol. 105, pp. 89–95, Jan. 2025.
- [9] R. Hernandez, M. Lopez, and A. Garcia, "Postural assessment using deep learning for physiotherapy applications," *Comput. Methods Programs Biomed.*, vol. 238, p. 107612, Feb. 2025.
- [10] Y. Cai, W. Li, and H. Zhang, "Swin-UNet for 3D human motion quality assessment in rehabilitation," *Med. Image Anal.*, vol. 89, p. 102876, Jan. 2025.
- [11] K. Rajesh, P. Kumar, and S. Sharma, "Spatio-temporal graph networks for exercise recognition using OpenPose," *Pattern Recognit. Lett.*, vol. 175, pp. 112–119, Nov. 2024.
- [12] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," in *Proc. NIPS Workshop Deep Learn.*, 2014.