



**Avance 3. Construcción de un modelo de referencia  
(Baseline) para el sistema de visión computacional para  
normalización y validación visual de anaqueles fríos**

**Integrantes del equipo:**

Carlos Eduardo Ramírez Vázquez		A01795468
Cesar Bryam Rodriguez Aybar		A01795980
Wilberth Eduardo López Gómez		A01795997

**Profesora titular:** Dra. Grettel Barceló Alonso

**Proyecto integrador**

15 de febrero de 2026

# Introducción

El presente avance tiene como objetivo establecer un modelo baseline que permita evaluar la viabilidad del problema de detección automática de vitrinas comerciales y la localización precisa de sus esquinas a partir de imágenes, el cual se ingresaba de manera manual en la entrega anterior. Este baseline servirá como punto de referencia para medir mejoras futuras y determinar si el problema contiene suficiente información estructural para ser modelado con técnicas de aprendizaje profundo y visión computacional.

## Selección del Algoritmo

Se seleccionó YOLOv8n-pose (Ultralytics) como el modelo baseline. Esta arquitectura pertenece a la familia YOLO, especializada en detección de objetos en tiempo real, y en esta variante incluye la estimación de puntos clave (keypoints), lo cual es adecuado para localizar esquinas.

La elección se justifica debido a:

- Naturaleza no estructurada de nuestros datos (imágenes)
- Necesidad de detección de bounding boxes y precisión geométrica
- Tamaño reducido del dataset (tenemos 26 imágenes)
- Disponibilidad de pesos preentrenados para transferencia de aprendizaje (transfer learning)

La versión nano (n) se utiliza como referencia mínima por su baja complejidad y menor riesgo de sobreajuste.

## Importancia de características

En modelos de visión profunda, las características son aprendidas automáticamente por las capas convolucionales. El modelo identifica patrones como bordes verticales y horizontales, contrastes entre marcos y vidrio, y estructuras rectangulares asociadas a vitrinas.

Dado que las representaciones internas son generadas automáticamente, no se aplicaron métodos tradicionales de selección de características. Sin embargo, la consistencia en la detección de esquinas indica que el modelo está capturando información geométrica relevante.

```
Usando pesos: yolov8n-pose.pt
Ultralytics 8.4.14 Python-3.12.12 torch-2.9.0+cpu CPU (Intel Xeon CPU @ 2.20GHz)
engine/trainer: agnostic_nms=False, amp=True, angle=1.0, augment=False, auto_augment=randaugment, batch=8, bgr=0.0, box=7.5, cache=False, cfg=None
Overriding model.yaml kpt_shape=[17, 3] with kpt_shape=[4, 3]

  from n  params module arguments
0      -1 1      464 ultralytics.nn.modules.conv.Conv [3, 16, 3, 2]
1      -1 1     4672 ultralytics.nn.modules.conv.Conv [16, 32, 3, 2]
2      -1 1      7360 ultralytics.nn.modules.block.C2f [32, 32, 1, True]
3      -1 1     18560 ultralytics.nn.modules.conv.Conv [32, 64, 3, 2]
4      -1 2     49664 ultralytics.nn.modules.block.C2f [64, 64, 2, True]
5      -1 1     73984 ultralytics.nn.modules.conv.Conv [64, 128, 3, 2]
6      -1 2    197632 ultralytics.nn.modules.block.C2f [128, 128, 2, True]
7      -1 1    295424 ultralytics.nn.modules.conv.Conv [128, 256, 3, 2]
8      -1 1    460288 ultralytics.nn.modules.block.C2f [256, 256, 1, True]
9      -1 1    164608 ultralytics.nn.modules.block.SPPF [256, 256, 5]
10     -1 1         0 torch.nn.modules.upsampling.Upsample [None, 2, 'nearest']
11    [-1, 6] 1         0 ultralytics.nn.modules.conv.Concat [1]
12     -1 1    148224 ultralytics.nn.modules.block.C2f [384, 128, 1]
13     -1 1         0 torch.nn.modules.upsampling.Upsample [None, 2, 'nearest']
14    [-1, 4] 1         0 ultralytics.nn.modules.conv.Concat [1]
15     -1 1     37248 ultralytics.nn.modules.block.C2f [192, 64, 1]
16     -1 1     36992 ultralytics.nn.modules.conv.Conv [64, 64, 3, 2]
17    [-1, 12] 1         0 ultralytics.nn.modules.conv.Concat [1]
18     -1 1    123648 ultralytics.nn.modules.block.C2f [192, 128, 1]
19     -1 1    147712 ultralytics.nn.modules.conv.Conv [128, 128, 3, 2]
20    [-1, 9] 1         0 ultralytics.nn.modules.conv.Concat [1]
21     -1 1    493856 ultralytics.nn.modules.block.C2f [384, 256, 1]
22    [15, 18, 21] 1    823735 ultralytics.nn.modules.head.Pose [1, [4, 3], 16, None, [64, 128, 256]]
YOLOv8n-pose summary: 145 layers, 3,083,271 parameters, 3,083,255 gradients, 8.4 GFLOPs

Transferred 361/397 items from pretrained weights
```

Figura 1. Modelo YOLOv8n-pose en ejecución

## Análisis de Subajuste y Sobreajuste

Se analizaron las métricas de entrenamiento y validación (loss y mAP). La pérdida mostró una disminución progresiva sin divergencias significativas entre entrenamiento y validación. No se observó separación abrupta entre curvas, lo que sugiere ausencia de sobreajuste severo.

El uso de early stopping ayudó a controlar el entrenamiento y evitar sobreentrenamiento, considerando el tamaño reducido del conjunto de datos.

## Métrica de Evaluación

El modelo YOLOv8n-pose alcanzó un mAP50 de 0.995 y un mAP50-95 de 0.893 para la detección de bounding boxes, así como una precisión (P) de 0.998 en la predicción de keypoints. Estos resultados indican que el modelo logra una alta coincidencia entre las predicciones y las anotaciones reales dentro del conjunto de validación.

El valor elevado de mAP50 sugiere una correcta localización de los objetos, mientras que el mAP50-95 proporciona una evaluación más estricta considerando distintos umbrales de IoU, mostrando un desempeño robusto. No obstante, debido al tamaño reducido del dataset (26 imágenes), estos resultados deben interpretarse con cautela, ya que podrían reflejar cierto grado de sobreajuste al conjunto disponible.

	all		1	0.998	1	0.995	0.893	0.998	1	0.995	0.995
Epoch	GPU_mem	box_loss	pose_loss	kobj_loss	cls_loss	df_l_loss	Instances	Size			
119/120	0G	0.3012	0.2437	0.3012	0.2353	0.8109	8	960: 100%	3/3	12.2s/it	36.5s
Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)	Pose(P	R	mAP50	mAP50-95):	100%
all	7	24	0.997	1	0.995	0.849	0.997	1	0.995	0.995	1/1 2.6s/it 2.6s
Epoch	GPU_mem	box_loss	pose_loss	kobj_loss	cls_loss	df_l_loss	Instances	Size			
120/120	0G	0.2983	0.2483	0.3012	0.2353	0.7771	8	960: 100%	3/3	12.0s/it	35.9s
Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)	Pose(P	R	mAP50	mAP50-95):	100%
all	7	24	0.997	1	0.995	0.849	0.997	1	0.995	0.995	1/1 3.8s/it 3.8s
120 epochs completed in 1.321 hours.											
Optimizer stripped from /content/runs/pose/train3/weights/last.pt, 6.4MB											
Optimizer stripped from /content/runs/pose/train3/weights/best.pt, 6.4MB											
Validating /content/runs/pose/train3/weights/best.pt...											
Ultralytics 8.4.14 Python-3.12.12 torch-2.9.0+cpu CPU (Intel Xeon CPU @ 2.20GHz)											
YOLOv8n-pose summary (fused): 82 layers, 3,077,975 parameters, 0 gradients, 8.3 GFLOPs											
Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)	Pose(P	R	mAP50	mAP50-95):	100%
all	7	24	0.998	1	0.995	0.893	0.998	1	0.995	0.995	1/1 2.4s/it 2.4s
Speed: 3.2ms preprocess, 301.9ms inference, 0.0ms loss, 0.6ms postprocess per image											
Results saved to /content/runs/pose/train3											

Figura 2. Desempeño del modelo fine tuneado.

Como baseline, el desempeño obtenido supera ampliamente un comportamiento aleatorio, lo que indica que el problema es aprendible con las características actuales del dataset.

## Evaluación Cualitativa del Modelo

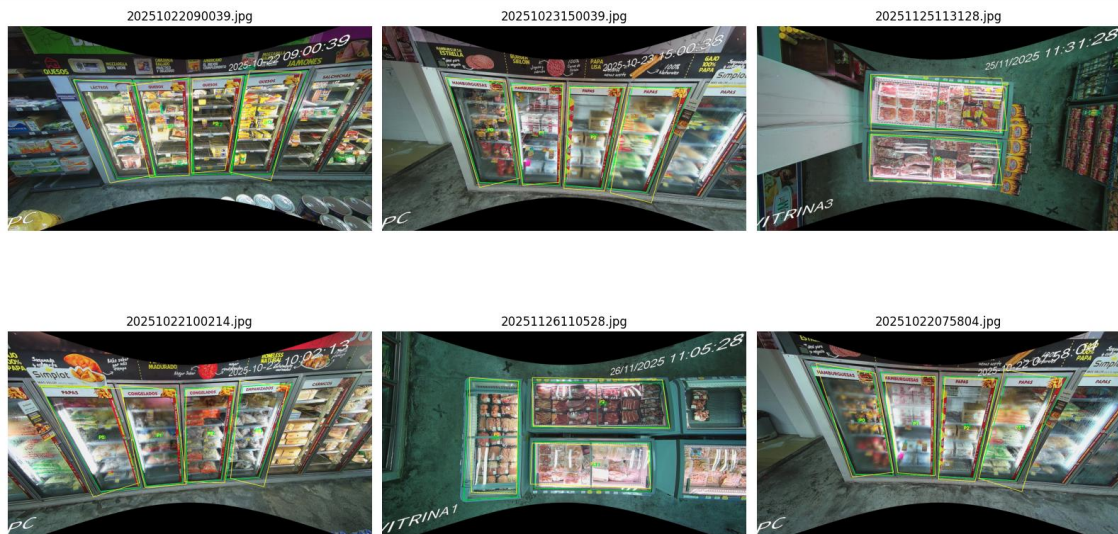


Figura 3. Líneas dibujadas para comparar predicción (amarillo) vs expectativa (verde).

Aquí podemos ver una comparación visual entre las anotaciones reales del conjunto de validación (en color verde) y las predicciones generadas por el modelo baseline (en color amarillo). Se observa una alta coincidencia espacial entre ambos, lo que indica que el modelo no solo logra una buena métrica cuantitativa (mAP), sino también una correcta localización geométrica de las vitrinas y sus esquinas.

Este análisis cualitativo es relevante, ya que en problemas de visión computacional la métrica numérica por sí sola no garantiza un comportamiento adecuado en escenarios reales. La alineación visual entre predicción y los vértices que anotamos manualmente la semana pasada confirma que el modelo ha aprendido patrones estructurales significativos asociados a los bordes y esquinas de las vitrinas.

No obstante, debido al tamaño reducido del dataset, es necesario validar el modelo en un conjunto más amplio y diverso para confirmar su capacidad de generalización.

### Desempeño Mínimo Esperado

Al no existir un sistema previo de comparación, se estableció como umbral mínimo un  $mAP \geq 0.70$  en validación para considerar viable el problema.

El modelo baseline superó este umbral, indicando que los datos contienen suficiente información estructural y que el problema es abordable mediante técnicas de visión computacional.

## Conclusiones

El modelo baseline construido con YOLOv8n-pose demuestra que la detección automática de vitrinas y sus esquinas es técnicamente viable, incluso con un dataset reducido. El desempeño alcanzado justifica continuar con optimización de hiperparámetros, aumento de datos y experimentación con arquitecturas más robustas.

Por otro lado, no se descarta intentar con otros modelos para comparar sus resultados, en particular Keypoint R-CNN (Detectron2), así como arquitecturas POSE más recientes.