

YOLOv12 Pill Counting: Attention-Based Object Detection Meets Pharmacy Automation

Kabui, Charles

2026-06-06

[Read at ToKnow.ai](#)

**YOLOv12 Pill Counting:
Attention-Based Detection
Meets Pharmacy Automation**

Attention-centric YOLO achieves 40.6% mAP at 1.64ms on T4 GPU

40.6% mAP on COCO beating YOLOv11-N by +1.2%	1.64ms Inference latency on T4 GPU (TensorRT)	1 in 30 Patients harmed by medication errors (WHO)
---	--	---

June 6, 2026

ToKnow.ai

YOLOv12, the latest in the YOLO (You Only Look Once) object detection series, replaces CNN-based feature extraction with an attention-centric architecture. Two key innovations make this work: an Area Attention mechanism that divides feature maps into segments for efficient parallel processing, and R-ELAN (Residual Efficient Layer Aggregation Network), which improves how the model fuses features across layers. The result is 40.6% mAP on

COCO with 1.64ms inference on a T4 GPU, beating YOLOv11-N by 1.2% mAP at comparable speed. The S-size variant beats RT-DETR-R18 while running 42% faster with only 36% of the computation. A hands-on tutorial from Labellerr demonstrates fine-tuning YOLOv12 for real-time pill counting, using the Ultralytics framework to train a custom model that detects and counts individual pills in images. [The full notebook, with a companion video walkthrough, runs in Google Colab with no local setup required.](#)

According to the WHO, medication-related harm affects 1 in every 30 patients, with over half of it preventable. Automated pill counting at the point of dispensing could catch errors before they reach patients. What makes this practical now is that YOLOv12's attention mechanism handles overlapping, similarly-colored small objects (exactly what pills are) better than its predecessors. The Labellerr tutorial takes this from research paper to working prototype: anyone with a Colab account can train a pill detection model in under an hour.

The broader pattern here matters. YOLO was CNN-only for a decade. YOLOv12 proves attention mechanisms can match CNN speed while improving accuracy, and was accepted at NeurIPS 2025. For a comparison of how vision-language models are pushing object detection further, see [NVIDIA LocateAnything](#), which takes a different approach with parallel box decoding.

Sources:

- [YOLOv12: Attention-Centric Real-Time Object Detectors \(arXiv:2502.12524\)](#)
- [YOLOv12 GitHub Repository](#)
- [Pill Counting Using YOLOv12 Notebook \(Labellerr\)](#)
- [WHO Patient Safety Fact Sheet: Medication Errors](#)
- [Labellerr YOLOv12 Evaluation Blog](#)

*Disclaimer: For information only. Accuracy or completeness not guaranteed. Illegal use prohibited. Not professional advice or solicitation. **Read more:** [/terms-of-service](#)*