# Real-Time Person Detection and Description Systems with Integrated Semantic Analysis

Mokksh Kapur

Dept. of Computer Science and Engineering

SRM University Sonipat

Haryana, India

Email: mokkshkapur@gmail.com

*Abstract*—**Advancements in surveillance technologies are pivotal in addressing the growing demands for real-time monitoring and actionable insights across industries such as security, public safety, and traffic management. This review explores existing research and technologies in real-time person detection, attribute analysis, and semantic description generation, highlighting their applications, limitations, and gaps. It further presents a novel system integrating YOLOv8 for efficient object detection, OpenCV for targeted attribute analysis, and Qwen2VL for rich vision-language understanding, coupled with Sentence-BERT for semantic search. This integrated approach addresses key challenges including latency, scalability in processing complex queries, and dynamic environment adaptability. We discuss the technical advantages, computational considerations, particularly for scaling VLMs, and challenges in complex scenarios like dense crowds. Future directions include integrating emotion recognition, multilingual capabilities, and more sophisticated predictive analytics to enhance system utility and proactive security measures.**

## I. INTRODUCTION

The increasing prevalence of surveillance systems has transformed public safety, urban management, and event monitoring. However, traditional systems often provide limited actionable insights, failing to meet the demands for real-time detection combined with descriptive analysis and contextual understanding. This limits their effectiveness in dynamic, complex scenarios where simply detecting a person is insufficient.

This review focuses on recent advancements in systems aiming to bridge this gap by integrating real-time person detection with description generation. We emphasize the convergence of deep learning object detection, computer vision attribute analysis, and sophisticated vision-language models (VLMs). The primary objective is to identify current trends, technological capabilities, inherent challenges (including computational scalability and real-world robustness), and research gaps. Based on this analysis, we present a scalable, modular system design that integrates state-of-the-art components to provide richer, searchable insights compared to traditional surveillance tools or standalone detection models.

## II. LITERATURE REVIEW

This review employs a systematic approach: literature search across major databases (IEEE Xplore, Springer, CVPR) using keywords like *YOLOv8*, *real-time object detection*, *vision-language models*, *semantic surveillance*, *attribute analysis*; inclusion criteria focusing on real-time capabilities, semantic feature extraction, and system integration; data extraction of methodologies, performance metrics (mAP, latency), and findings; comparative analysis evaluating detection accuracy, processing speed, scalability, and adaptability; and synthesis to identify trends, challenges, and research gaps guiding the proposed system.

### A. AI-Driven Surveillance Platforms: Applications and Ethical Context

Large-scale AI surveillance deployments provide context for capabilities and societal considerations:

- **Clearview AI**: Demonstrates the power of large-scale facial recognition for law enforcement but faces intense scrutiny regarding privacy violations due to its controversial, unconsented scraping of public web images.
- **Amazon Rekognition & Microsoft Azure Cognitive Services**: Offer broad, cloud-based AI capabilities (detection, facial analysis) enabling scalable solutions. However, documented concerns regarding potential demographic biases in their algorithms and the potential for misuse highlight the critical need for fairness, accountability, and transparency in deploying such technologies at scale [24].

These examples underscore the demand for automated analysis but also establish a crucial backdrop of ethical responsibilities concerning data privacy and algorithmic bias that must be considered in developing and deploying any surveillance technology, including the system proposed herein.

### B. Core Technologies in Person Analysis

Advancing beyond basic detection requires integrating specialized technologies:

*1) Object Detection Frameworks: The Case for YOLOv8:* Real-time object detection is the foundation. While two-stage detectors like Faster R-CNN [4] improved accuracy via Region Proposal Networks (RPNs), their sequential nature often limits speed. YOLOv8 [1], representing the evolution of the You Only Look Once paradigm, offers distinct advantages for real-time applications. Its **single-stage architecture** processes the entire image grid simultaneously to predict bounding boxes and class probabilities directly, eliminating the separate

proposal generation step. Furthermore, YOLOv8 incorporates highly optimized **backbone networks (e.g., based on CSPNet principles)** and **feature fusion mechanisms (like PANet necks)**, along with techniques such as **anchor-free detection heads**. This architectural efficiency translates directly to significantly **higher inference speeds (Frames Per Second - FPS)** at comparable accuracy levels (mAP) on standard benchmarks [10] compared to most two-stage detectors. While transformer-based detectors like DETR [25] offer novel end-to-end approaches, their computational demands, particularly the self-attention mechanism's quadratic complexity, often result in higher latency, making highly optimized CNN architectures like YOLOv8 generally more suitable for resource-constrained edge deployments or applications with strict real-time demands.

*2) Attribute Analysis Techniques: Leveraging OpenCV:* Understanding attributes (gender, clothing type/color) adds valuable semantic detail. Lightweight models available via libraries like OpenCV [5], such as the widely used *face.caffemodel* and *gender_deploy.prototxt* (often based on foundational work like [8]), provide efficient means for basic attribute classification on detected regions. Their speed makes them suitable for integration into real-time pipelines. However, it is crucial to acknowledge their limitations in uncontrolled environments. These models often learn representations based on **holistic facial appearance or the expected geometric arrangement of key features**. Consequently, their reliability can decrease under challenging conditions common in dynamic surveillance:

- **Occlusion** (e.g., masks, sunglasses, hats, partial view due to crowds) directly obscures the features the model relies upon.
- **Significant Pose Variations** (non-frontal views) alter the appearance and relative position of facial features compared to the training data distribution.
- **Low Resolution and Poor Lighting** can degrade feature quality beyond the model's ability to discern relevant patterns.

Therefore, while useful for providing quick attribute hints, outputs from such lightweight models should be interpreted with caution, potentially requiring confidence thresholding or fusion with other cues in a robust system design.

*3) Vision-Language Models for Semantic Understanding:* Generating human-like descriptions requires deeper comprehension. Vision-Language Models (VLMs) excel here. Models like CLIP [26] learn aligned image-text embeddings useful for zero-shot tasks. More advanced VLMs like Qwen2VL [6] are specifically architected for tasks like Visual Question Answering (VQA) and generating detailed, context-aware image captions. Applied to detected person ROIs (as shown in Fig. 2, Fig. 3), these models can produce rich textual descriptions capturing appearance, inferred actions, and interactions, moving significantly beyond simple object tags or basic attributes [11], [12], [13], [14], [15].

## C. Synthesis and Identified Gaps

The review highlights potent individual components but reveals critical gaps in integration and deployment:

- **Integrated Systems:** Few systems seamlessly combine high-speed detection, targeted attribute analysis, *and* rich VLM-generated descriptions with semantic search capabilities into a cohesive, real-time pipeline.
- **Semantic Queryability:** Enabling intuitive, natural language search based on appearance, actions, or context described textually remains underdeveloped compared to metadata tag searching.
- **Robustness in the Wild:** Ensuring consistent performance of all components (especially attribute analysis and description) under real-world dynamism (lighting, pose, occlusion) is an ongoing challenge.
- **Efficiency and Scalability:** Balancing the high computational cost of advanced VLMs (significant **GPU memory demands and teraFLOP computations**) with real-time constraints is a major hurdle, especially for large-scale, multi-camera deployments.

These gaps underscore the need for architectures like the proposed system, focusing on modular integration, efficient component selection where possible, and strategies to manage computational load for advanced features. Table I provides further context by comparing related research efforts.

## III. TRENDS, CHALLENGES, AND RESEARCH GAPS

### A. Trends

- Increasing adoption of modular, scalable architectures.
- Integration of VLMs for enhanced descriptive analytics and semantic search.
- Dominance of high-performance real-time detectors like YOLO variants.

### B. Challenges

- Adapting systems robustly to dynamic environmental variations (lighting, occlusion, pose).
- Computational bottlenecks in attribute analysis pipelines and especially VLM inference at scale.
- Limited integration of advanced features (e.g., emotion recognition, multilingual support) in deployable systems.
- Handling high-density crowd scenarios effectively (discussed further in Section V).

### C. Research Gaps

- Seamless integration of real-time detection with deep semantic analysis and description.
- Development of intuitive, efficient query mechanisms for VLM-generated descriptive data.
- Scalable deployment strategies for computationally intensive components like VLMs across many streams.

## IV. PROPOSED SYSTEM

Addressing the identified gaps, the proposed system integrates YOLOv8 [1], OpenCV [5], Qwen2VL [6], and Sentence-BERT [27] into a cohesive pipeline for enhanced real-time person analysis (Fig. 1).

**Architecture Overview:** Input video frames are processed sequentially and concurrently. YOLOv8 performs initial person detection. For each detected person's ROI, lightweight attribute analysis (e.g., gender) is performed using OpenCV modules, while the more computationally intensive Qwen2VL model generates a detailed textual description. These outputs are associated with the detection instance. The generated descriptions and associated metadata are stored persistently. The descriptions are then embedded using Sentence-BERT, with the resulting embeddings indexed in a suitable data store for efficient semantic retrieval based on natural language queries.

**Component Roles:**

- **YOLOv8:** Provides the foundational high-speed, accurate person detection.
- **OpenCV:** Offers efficient extraction of specific, predefined attributes using fast, lightweight models.
- **Qwen2VL:** Generates rich, nuanced, context-aware textual descriptions, capturing details beyond simple attributes. Acknowledging its computational cost, its invocation might be triggered selectively or at intervals (see Section V).
- **Sentence-BERT:** Enables powerful semantic search by converting textual descriptions and user queries into a common vector space for similarity matching.

**Workflow:** The process involves:

1) Frame Ingestion.
2) YOLOv8 Detection to obtain bounding boxes.
3) ROI Extraction based on bounding boxes.
4) Parallel processing for each detected person:
    a. Pass ROI to OpenCV modules for attribute classification.
    b. Pass ROI to Qwen2VL to generate a textual description.
5) **Data Persistence:** Store the associated data (bounding box, timestamp, attributes, generated description) persistently.
6) **Indexing for Search:** Index the generated description using Sentence-BERT, storing the resulting vector embedding in a data store optimized for efficient similarity search (e.g., a vector database).
7) User provides a natural language query.
8) Embed the query using Sentence-BERT.
9) Perform vector similarity search against the indexed description embeddings in the data store.
10) Retrieve ranked results (matching persons with their associated data/frames).

This integrated workflow aims to deliver not just detection but actionable, searchable semantic insights managed through appropriate data persistence and indexing strategies.
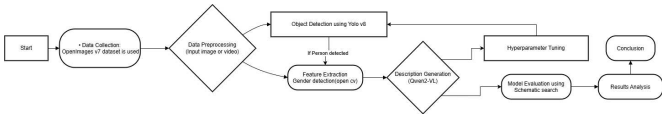


Fig. 1. System Architecture Flow Diagram: Illustrates the pipeline from input processing via YOLOv8 (Detection), OpenCV (Attribute Analysis), and Qwen2VL (Description Generation) to Semantic Search using Sentence-BERT, implying data storage/indexing.

## V. SYSTEM CAPABILITIES, CHALLENGES, AND COMPARATIVE POSITIONING

The proposed system architecture offers unique strengths but also faces practical challenges inherent in real-time, complex visual analysis.

**Capabilities:**

- **Integrated Semantic Analysis:** Moves beyond simple detection by combining attributes and rich VLM descriptions within a single pipeline.
- **Real-Time Feasibility:** Leverages the speed of YOLOv8 and allows strategies (e.g., frame analysis intervals, tested at 5s) to manage the computational load of components like Qwen2VL, enabling near real-time operation on capable hardware (preliminarily tested on Intel i7/RTX 4060).
- **Semantic Querying:** Integrates Sentence-BERT for intuitive natural language search over generated descriptions, enhancing data retrieval beyond keyword matching.

**Challenges and Considerations:**

- **VLM Scalability:** While feasible for single streams with interval processing, scaling the Qwen2VL description generation across numerous cameras concurrently presents a major computational challenge due to its high GPU/memory demands. Distributed processing, model optimization (quantization, distillation), or adaptive triggering mechanisms would be crucial for large-scale deployment.
- **High-Density Crowds:** Such scenarios pose difficulties. Severe inter-person occlusion can degrade YOLOv8's detection accuracy and makes isolating individuals for reliable attribute analysis (via OpenCV or VLM) extremely challenging. Furthermore, the computational load of generating descriptions via Qwen2VL increases linearly (or worse, if analysis complexity grows) with the number of detected persons, potentially creating bottlenecks in dense scenes. Strategies like selective description generation (e.g., focusing on individuals near an event zone or matching a basic query filter) might be necessary.
- **Robustness of Components:** As discussed (Sec II.B.2), lightweight attribute analyzers have inherent limitations. While VLMs are generally more robust, their descriptive accuracy can also degrade under extreme occlusion, very low resolution, or unusual contexts not well-represented in their training data. System reliability depends on the robustness of each pipeline stage.

**Comparative Positioning:** This system distinguishes itself through its specific focus on integrating multiple analysis modalities for enriched, searchable output. Compared to **cloud platforms** (Rekognition, Azure), it prioritizes potential edge deployment (reducing latency) and offers deeper semantic descriptions via VLM coupled with dedicated semantic search, rather than just predefined tags/attributes. Unlike **specialized systems** (Clearview AI) focused narrowly on facial identification, our approach targets broader person analysis (appearance, context) without relying on potentially privacy-invasive identity databases. Relative to **standard detection systems**, it adds significant value through integrated attribute analysis, detailed VLM descriptions, and semantic retrieval. While research studies in Table I explore individual aspects, this work emphasizes the **synergistic integration** of these components into a functional pipeline designed for enhanced situational awareness. Our preliminary evaluation confirms the *feasibility* of this integration, while acknowledging that rigorous, quantitative benchmarking against these diverse systems requires standardized protocols and further development.

## VI. CONCLUSION AND FUTURE DIRECTIONS

This paper reviewed real-time person detection and description systems, highlighting advancements and persistent gaps, particularly in integrating deep semantic analysis with efficient detection and querying. We proposed a system architecture combining YOLOv8, OpenCV, Qwen2VL, and Sentence-BERT, offering a pathway towards more insightful and searchable surveillance analytics. The system demonstrates the feasibility of integrating these powerful components, providing capabilities beyond traditional detection, including attribute analysis, rich textual descriptions, and natural language semantic search. Key strengths lie in this integration, its adaptability through modular design, and the potential for real-time operation with appropriate hardware and processing strategies.

Acknowledging the computational challenges, especially in scaling VLMs and handling dense crowds, and the inherent limitations of individual components in extreme conditions, significant avenues for future work exist:

- **Predictive Analytics:** Moving beyond reactive analysis by integrating temporal reasoning. Future work includes leveraging aggregated data (e.g., trajectories derived from tracking detections over time, recurring semantic descriptions of behavior/appearance) to train models capable of **forecasting potential security events or anomalies**. This could involve identifying unusual loitering patterns, detecting anomalous crowd dynamics, or recognizing sequences of described actions potentially indicative of risks, enabling proactive intervention strategies.
- **Emotion Recognition:** Integrating models to analyze facial expressions or body language for more nuanced behavioral understanding and context-aware insights.
- **Multilingual Support:** Adapting VLM components and query interfaces to support multiple languages, broadening accessibility and applicability.

- **Optimization for Scalability:** Developing and evaluating techniques like model quantization, knowledge distillation, efficient multi-stream scheduling, and adaptive triggering for VLM inference to make large-scale deployment more practical.
- **Enhanced Robustness:** Investigating methods to improve component robustness in challenging conditions, potentially through multi-modal fusion or uncertainty-aware processing.
- **Rigorous Evaluation:** Conducting comprehensive benchmarking on standardized datasets to quantify end-to-end performance, including detection accuracy, description quality (BLEU, ROUGE), attribute accuracy, search precision/recall, and latency under various loads.

This work lays a foundation for developing more intelligent surveillance systems capable of deeper understanding and interaction, while carefully considering the associated technical and ethical dimensions.



Fig. 2. "The person in the image is wearing a patterned dress. The dress appears to be of a dark color with a floral or paisley design. The dress covers the person's upper body and legs." (Example Qwen2VL output)
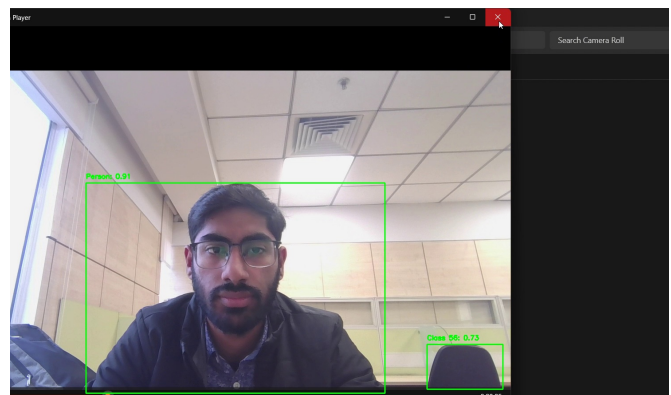


Fig. 3. "The person in the image is wearing glasses and a dark-colored shirt. The shirt appears to be a solid color, possibly black or dark blue. The glasses are a dark frame, which contrasts with the lighter background of the shirt. The person's hair is short and neatly styled. The overall appearance is professional and neat." (Example Qwen2VL output)

REFERENCES

[1] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," arXiv preprint arXiv:1804.02767, 2018.

[2] A. Vaswani, et al., "Attention is All You Need," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[3] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2015.

[4] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[5] OpenCV Documentation, [Online]. Available: https://opencv.org.

[6] Qwen Team, "Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and More," arXiv preprint arXiv:2308.12966, 2023.

[7] *Consider citing [4] again for RPN concepts or find specific FCN paper if distinct concept intended.*

[8] R. Rothe, R. Timofte, L. Van Gool, "Deep expectation of real and apparent age from a single image without facial landmarks," *International Journal of Computer Vision (IJCV)*, vol. 126, pp. 144–157, 2018.

[9] M. Sharma and P. Goswami, "Efficient Support System for Hearing Disabled Using CNN," in *Proc. IEEE Int. Conf. on Advances in Computing, Communication Control and Networking (ICAC3N)*, 2023.

[10] P. Goswami, et al., "Real-time Evaluation of Object Detection Models," *Applied Soft Computing*, vol. 152, p. 111234, 2024.

[11] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and Tell: A Neural Image Caption Generator," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3156–3164.

[12] A. Karpathy, A. Joulin, and L. Fei-Fei, "Deep Fragment Embeddings for Bidirectional Image Sentence Mapping," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2014.

[13] R. Bernardi, et al., "Automatic Description Generation from Images: A Survey of Models, Datasets, and Evaluation Measures," *J. Artif. Intell. Res. (JAIR)*, vol. 55, pp. 409–442, 2016.

[14] M. Mitchell, et al., "Midge: Generating Image Descriptions from Computer Vision Detections," in *Proc. 13th Conf. of the European Chapter of the Association for Computational Linguistics (EACL)*, 2012, pp. 747–756.

[15] G. Kulkarni, et al., "Baby Talk: Understanding and Generating Simple Image Descriptions," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 12, pp. 2891–2903, Dec. 2013.

[16] M. Ahmad, I. Ahmed, K. Ullah, I. Khan, A. Khattak, and A. Adnan, "Person Detection from Overhead View: A Survey," *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 4, 2019.

[17] J. Cao, L. Sun, M. G. Odoom, F. Luan, and X. Song, "Counting People by Using a Single Camera without Calibration," in *Proc. Chinese Control and Decision Conference (CCDC)*, IEEE, 2016, pp. 2048–2051.

[18] J. Nalepa, J. Szymanek, and M. Kawulok, "Real-time People Counting from Depth Images," in *Proc. Int. Conf. Beyond Databases, Architectures and Structures (BDAS)*, Springer, 2015, pp. 387–397.

[19] L. Del Pizzo, et al., "Counting People by RGB or Depth Overhead Cameras," *Pattern Recognition Letters*, vol. 81, pp. 41–50, 2016.

[20] V. Carletti, et al., "An Efficient and Effective Method for People Detection from Top-View Depth Cameras," in *Proc. IEEE Int. Conf. on Advanced Video and Signal-Based Surveillance (AVSS)*, 2017, pp. 1–6.

[21] *Consider citing specific benchmark datasets like MOTChallenge, CrowdHuman or specific counting benchmark papers if applicable.* H. Lin, et al., "Cross-dataset People Counting using Synthetic Data," *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, could be relevant.

[22] H. Zhang, et al., "AI-driven behavior biometrics framework for robust human activity recognition in surveillance systems," *Engineering Applications of Artificial Intelligence*, vol. 127, Part A, p. 107218, 2024. [Online]. Available: https://doi.org/10.1016/j.engappai.2023.107218.

[23] A. Singh, et al., "Human action recognition using attention-based LSTM network with dilated CNN features," *Future Generation Computer Systems*, vol. 125, pp. 820-830, 2021. [Online]. Available: https://doi.org/10.1016/j.future.2021.06.045.

[24] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in *Proc. Conf. on Fairness, Accountability, and Transparency (FAT*)*, 2018.

[25] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, "End-to-End Object Detection with Transformers," in *Proc. European Conf. on Computer Vision (ECCV)*, 2020.

[26] A. Radford, et al., "Learning Transferable Visual Models From Natural Language Supervision," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2021.

[27] N. Reimers, I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," in *Proc. Conf. on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

TABLE I
COMPARATIVE ANALYSIS OF SELECTED STUDIES

| Study | Objective | Methodology | Advantages | Limitations |
|---|---|---|---|---|
| Efficient Support System for Hearing Disabled Using CNN [9] | Gesture detection | CNN-based real-time detection | High accuracy (99%), real-time processing | Limited to static gestures, dataset-specific |
| Real-time Evaluation of Object Detection Models [10] | Compare YOLOv8, Faster R-CNN, DETR | mAP, inference speed comparison | Detailed performance metrics on diverse datasets | DETR struggles with large datasets, slower inference |
| Automatic Description Generation from Images [13] | Generate textual descriptions | Vision-language models | High-quality captions, semantic search potential | Often limited to offline processing, lacks real-time adaptability |
| Age and Gender Prediction Using Caffe and OpenCV [8] | Predict age/gender | Caffe model and OpenCV | Lightweight and efficient | Sensitive to low-res images, pose, occlusion (as discussed) |
| Fully Convolutional Region Proposal Networks (Concept Ref [4]) | Enhance region proposals | Convolutional networks for proposal generation | Improves Faster R-CNN efficiency | Adds computational overhead compared to single-stage |
| Deep Learning Object Detection Surveys (e.g., [16]) | Analyze advancements | Comparison of architectures | Comprehensive overview | Often lack specific real-time integration/implementation insights |
| Proposed Solution (This Work) | Real-time detection, attribute analysis, description, semantic search | Integrated YOLOv8, OpenCV, Qwen2VL, Sentence-BERT | Real-time feasibility, integrated semantic richness, queryability | Requires tuning, GPU needed (VLM), preliminary eval, known challenges (crowds, VLM scaling) |

TABLE II
MODEL COMPARISON (YOLOV8 VS ALTERNATIVES)

| Model/Feature | Primary Task | Type | Real-time Perf. | Accuracy Metric (Typical) | Limitations |
|---|---|---|---|---|---|
| **YOLOv8** [1] | Object Detection | Real-time detector (1-stage CNN) | Very Fast (low latency) | mAP (detection) High | Struggles w/ extreme occlusion/overlap; Less context than VLM |
| **Faster R-CNN** [4] | Object Detection | 2-stage detector w/ RPN | Slower than YOLO | mAP (detection) High | Higher computational cost, slower inference latency |
| **DETR** [25] | Object Detection | Transformer-based detector | Slowest | mAP (detection) Competitive | Higher latency, complex training, struggles dense scenes |
| **ViTs (general)** [26] | Classification / Feature Extraction | Transformer-based vision model | Moderate | Accuracy (classification) High | Needs large datasets; Often used as backbone, not standalone detector |
| **Qwen2VL** [6] | VQA / Description Gen | Vision-Language Model | Slow | BLEU/ROUGE/METEOR (text gen) High | High computational cost (GPU required), scaling challenge, potential biases |

Note: Performance and accuracy metrics are task-dependent. Direct comparison across different primary tasks requires careful consideration of relevant benchmarks. Scalability and dynamic environment flexibility vary based on specific implementation and fine-tuning. ViTs are often used as backbones (e.g., within CLIP or some detectors).