

ArmFormer: Lightweight Transformer Architecture for Real-Time Multi-Class Weapon Segmentation and Classification

Akhila Kambhatla^{*,†,1,2}, Taminul Islam^{†,1,2}, and Khaled R Ahmed^{1,2}

¹School of Computing, Southern Illinois University, Carbondale, IL, USA

²BASE Lab, Southern Illinois University, Carbondale, IL, USA
{akhila.kambhatla, taminul.islam, khaled.ahmed}@siu.edu

Abstract—The escalating threat of weapon-related violence necessitates automated detection systems capable of pixel-level precision for accurate threat assessment in real-time security applications. Traditional weapon detection approaches rely on object detection frameworks that provide only coarse bounding box localizations, lacking the fine-grained segmentation required for comprehensive threat analysis. Furthermore, existing semantic segmentation models either sacrifice accuracy for computational efficiency or require excessive computational resources incompatible with edge deployment scenarios. This paper presents ArmFormer, a lightweight transformer-based semantic segmentation framework that strategically integrates Convolutional Block Attention Module (CBAM) with MixVisionTransformer architecture to achieve superior accuracy while maintaining computational efficiency suitable for resource-constrained edge devices. Our approach combines CBAM-enhanced encoder backbone with attention-integrated hamburger decoder to enable multi-class weapon segmentation across five categories: handgun, rifle, knife, revolver, and human. Comprehensive experiments demonstrate that ArmFormer achieves state-of-the-art performance with 80.64% mIoU and 89.13% mFscore while maintaining real-time inference at 82.26 FPS. With only 4.886G FLOPs and 3.66M parameters, ArmFormer outperforms heavyweight models requiring up to 48× more computation, establishing it as the optimal solution for deployment on portable security cameras, surveillance drones, and embedded AI accelerators in distributed security infrastructure.

Index Terms—weapon detection, semantic segmentation, transformer, attention mechanism, real-time processing, computer vision

I. INTRODUCTION

The escalating threat of firearm violence poses a critical challenge to public safety, with over 42,155 firearm-related deaths reported in the United States in 2024-2025 [1]. Traditional security infrastructure relying on human surveillance and manual threat assessment has proven inadequate for addressing modern security challenges, particularly in high-traffic environments such as airports, educational institutions, and urban centers where potential threats exceed human operational capabilities. Surveillance systems require continuous human monitoring, yet concentration deteriorates to 83%, 84%, and 64% after one hour when monitoring 4, 9, and 16

screens respectively [2], underscoring the need for automated threat detection systems.

Current weapon detection methods predominantly employ object detection frameworks providing bounding box localizations [3], [4], yet fail to deliver pixel-level precision necessary for accurate threat assessment in complex scenarios involving occlusions [5], varying illumination [6], and multiple weapon types [7], [8]. While YOLO-based approaches demonstrate effectiveness in general object detection [9], they struggle with small and partially occluded objects [3], [10], which are common in weapon detection scenarios. Semantic segmentation addresses these limitations by providing pixel-level classification and precise boundary delineation essential for accurate threat assessment [11], [12].

Recent transformer-based semantic segmentation architectures such as SegFormer [13] have achieved superior performance in general computer vision tasks, demonstrating state-of-the-art results on benchmarks like ADE20K and Cityscapes with 5× speed improvements [14]. However, these advances remain largely unexplored for security-critical applications requiring specialized multi-class weapon identification, with current research predominantly relying on traditional CNN-based approaches [15]–[17]. The absence of transformer-based methodologies in weapon detection is particularly notable given their demonstrated effectiveness in modeling complex spatial relationships and extracting multi-resolution features [18], suggesting significant untapped potential for enhancing both accuracy and operational efficiency.

Despite revolutionizing semantic segmentation tasks, existing transformer architectures exhibit fundamental deficiencies when applied to weapon detection: excessive computational overhead incompatible with real-time requirements [19]–[21], lack of specialized attention mechanisms for weapon classification, and performance degradation when distinguishing morphologically similar weapon categories under challenging environmental conditions [5]–[7]. Multi-class weapon classification is essential for threat assessment, as distinguishing weapon types (handgun, rifle, knife, revolver) enables security systems to evaluate threat levels, operational ranges, and deploy proportionate countermeasures—whereas

*Corresponding author.

binary classification offers insufficient information for nuanced decision-making in critical security situations. In response to these limitations, this research introduces ArmFormer. This transformer-based segmentation architecture employs strategic Convolutional Block Attention Module (CBAM) [22] integration across encoder-decoder pathways to enhance multi-class weapon recognition while preserving computational tractability. Our methodology leverages Google Open Images and IMFDB database [23] to enable granular pixel-level classification of weapon classes and human subjects. Experimental validation demonstrates that ArmFormer achieves superior performance compared to established baseline models while maintaining the inference efficiency required for real-world security deployment.

II. RELATED WORK

Weapon detection methodologies have evolved from traditional sensor-based approaches to sophisticated deep learning architectures, yet achieving pixel-level precision with real-time efficiency for multi-class weapon segmentation remains an open challenge.

A. Weapon Detection: From Traditional to Deep Learning Approaches

Early weapon detection systems employed traditional techniques including radar [24], millimeter wave imaging [25], and thermal/infrared sensors [26] for security screening in airports and public venues, with TSA intercepting 6,678 firearms at checkpoints in 2024 [27]. These methods utilized pattern matching [28], density descriptors [29], and cascade classifiers [30], but suffered from high false positive rates due to object orientation variability and required continuous human monitoring of surveillance feeds [31].

The advent of deep learning transformed weapon detection through automated feature extraction. Classical computer vision approaches using hand-crafted descriptors such as HOG [32] and SIFT [33] with conventional classifiers were succeeded by CNN-based architectures including VGGNet [34] and ResNet [35] for weapon recognition [36]. Single-stage object detection frameworks, particularly YOLO variants (YOLOv3, YOLOv4, YOLOv8), have demonstrated superior real-time performance [37]–[39]. However, these methods present inherent limitations: coarse bounding box localization lacking pixel-level granularity, diminished accuracy on small-scale weapons, and insufficient robustness against partial occlusions [40]. Furthermore, CNNs’ limited receptive fields restrict their capacity to model long-range spatial dependencies essential for comprehensive scene understanding [41], [42], while transformers address these limitations through self-attention mechanisms enabling global spatial relationship modeling [43].

B. Transformer-based Semantic Segmentation and Attention Mechanisms

Vision Transformers (ViTs) [21] have revolutionized dense prediction tasks by establishing self-attention as an effective alternative to convolutions for pixel-level classification.

Semantic segmentation has progressed from FCNs [44] and U-Net [45] to transformer-based frameworks like SegFormer [13], Segmenter [46], Swin [20], achieving superior performance through hierarchical feature representation with efficient mix-FFN designs and overlapped patch embeddings [47]. While transformer-based approaches have shown promise for weapon detection through DETR-based localization frameworks [48], their application to pixel-level weapon segmentation remains largely unexplored, with computational efficiency challenges limiting deployment on resource-constrained security infrastructure.

Attention mechanisms have evolved from Squeeze-and-Excitation (SE) blocks to sophisticated spatial-channel integration schemes enhancing feature discriminability [49], [50]. The Convolutional Block Attention Module (CBAM) represents a significant advancement, systematically combining channel and spatial attention pathways for sequential feature refinement, demonstrating enhanced selectivity and superior multi-scale object handling [51]–[53]. Recent weapon detection research has shifted toward targeted attention mechanisms rather than complete transformer architectures, as full self-attention computations prove too expensive for real-time security systems. However, attention mechanism exploration remains limited in security-critical applications, with most research focusing on general-purpose vision tasks rather than specialized weapon detection requirements.

C. Research Gap and Motivation

Current transformer-based models exhibit critical deficiencies for weapon detection: excessive computational overhead incompatible with real-time deployment, lack of domain-specific optimizations for distinguishing morphologically similar weapon categories, and suboptimal attention mechanisms designed for general vision rather than security-critical applications [5]–[7], [19]–[21]. The fundamental challenge lies in balancing pixel-level precision with computational efficiency, as existing models sacrifice inference speed for accuracy. No prior work has combined transformers with pixel-level semantic segmentation specifically for multi-class weapon classification, despite transformers’ proven effectiveness in general segmentation tasks [13], [46] and weapon detection exploration through CNN-based [15], [16] or DETR-based [54] approaches. ArmFormer addresses this gap by delivering precise pixel-level weapon classification while maintaining computational tractability essential for operational security systems.

III. DATASET AND PRE-PROCESSING

A lack of a consistent dataset for firearm segmentation and recognition prompted the development of a unique dataset consisting of 8097 images curated from Google Open Images and the IMFDB [23] database. High-quality weapon images with diverse perspectives were carefully selected to ensure effective real-world detection performance. Preprocessing involved systematic removal of background noise and irrelevant elements from each image using computational tools to

TABLE I
DATASET SAMPLES DISTRIBUTION (COUNT)

Class	Test	Train	Validation	Total
Handgun	215	1506	440	2161
Human	133	928	269	1330
Knife	145	1017	288	1450
Revolver	172	1280	345	1797
Rifle	167	1115	328	1610

enhance dataset quality and training relevance. The dataset annotation process employs a sophisticated semi-automated approach utilizing the Segment Anything Model 2 (SAM2) framework [55] for multi-class object segmentation and has demonstrated robust multi-class segmentation performance across diverse domains [56]. The pipeline standardizes input image to 640×640 and implements an interactive annotation strategy that combines minimal human supervision with automated propagation across image sequences. SAM2’s video/image prediction capabilities propagate the initial annotations across the entire image sequence, maintaining object identity and segmentation consistency throughout temporal sequences. The model generates both individual class-specific masks and combined multi-class masks. Each class is assigned unique identifiers and corresponding grayscale values for mask generation (*Handgun: 51, Human: 102, Knife: 153, Rifle: 204, Revolver: 255*), enabling precise class differentiation in the ground truth annotations while maintaining computational efficiency. Table [I] presents the distribution of samples across training, testing, and validation sets for each weapon class, while Figure 1 illustrates representative examples from each class alongside their corresponding ground truth masks and segmented outputs.

IV. METHODOLOGY

A. ArmFormer Architecture

We propose ArmFormer, a lightweight transformer-based semantic segmentation framework that strategically integrates Convolutional Block Attention Module (CBAM) [22] with MixVisionTransformer [57] architecture to achieve superior feature representation capabilities for real-time multi-class weapon detection. As illustrated in Figure 2, the architecture employs a hierarchical encoder-decoder paradigm with attention mechanisms integrated at both encoder and decoder stages to enhance feature discrimination while maintaining computational efficiency suitable for edge deployment.

1) *CBAM Attention Mechanism*: CBAM serves as the core attention mechanism in ArmFormer, enhancing feature representation through sequential channel and spatial attention refinement. Given an input feature map $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ with C channels, height H , and width W , CBAM first applies channel attention to identify important feature channels, followed by spatial attention to localize informative regions.



Fig. 1: Visualization of original images, masked annotations, and segmented results for each weapon category including handgun, human, knife, revolver, and rifle.

The channel attention module exploits inter-channel relationships by aggregating spatial information through both average and max pooling operations. The channel attention map $\mathbf{M}_c \in \mathbb{R}^{C \times 1 \times 1}$ is computed as:

$$\mathbf{M}_c = \sigma(\text{MLP}(\text{AvgPool}(\mathbf{F})) + \text{MLP}(\text{MaxPool}(\mathbf{F}))) \quad (1)$$

where σ denotes sigmoid activation and MLP is a shared multi-layer perceptron with reduction ratio $r = 16$ that compresses and expands channel dimensions. The channel-refined feature is obtained as $\mathbf{F}' = \mathbf{M}_c \otimes \mathbf{F}$.

Following channel attention, the spatial attention module computes attention weights across spatial locations. The spatial attention map $\mathbf{M}_s \in \mathbb{R}^{1 \times H \times W}$ is generated by:

$$\mathbf{M}_s = \sigma(\text{Conv}^{7 \times 7}([\text{AvgPool}_c(\mathbf{F}'); \text{MaxPool}_c(\mathbf{F}')])) \quad (2)$$

where pooling operates along the channel dimension, and the concatenated descriptors are processed through a 7×7 convolution with kernel size $k = 7$. The final CBAM-refined feature output is:

$$\mathbf{F}_{out} = \mathbf{M}_s \otimes \mathbf{F}' = \mathbf{M}_s \otimes (\mathbf{M}_c \otimes \mathbf{F}) \quad (3)$$

This sequential attention mechanism allows the network to adaptively emphasize discriminative features while suppressing irrelevant background information, crucial for accurate weapon detection in complex security scenarios.

2) *CBAM-Enhanced Encoder Backbone*: The encoder backbone implements a CBAM-enhanced MixVisionTransformer architecture consisting of four hierarchical stages with progressive spatial resolution reduction and channel expansion to capture multi-scale contextual features. The input image $\mathbf{I} \in \mathbb{R}^{3 \times H_0 \times W_0}$ is processed through hierarchical stages with the following channel configurations: Stage 1 produces features with 32 channels, Stage 2 expands to 64 channels, Stage 3 increases to 160 channels, and Stage 4 reaches 256 channels, corresponding to spatial resolutions of $H_0/4$, $H_0/8$, $H_0/16$, and $H_0/32$ respectively.

Each stage employs overlapped patch embeddings with stride 4 and kernel size 7 to partition the feature map into overlapping patches, preserving local continuity while enabling

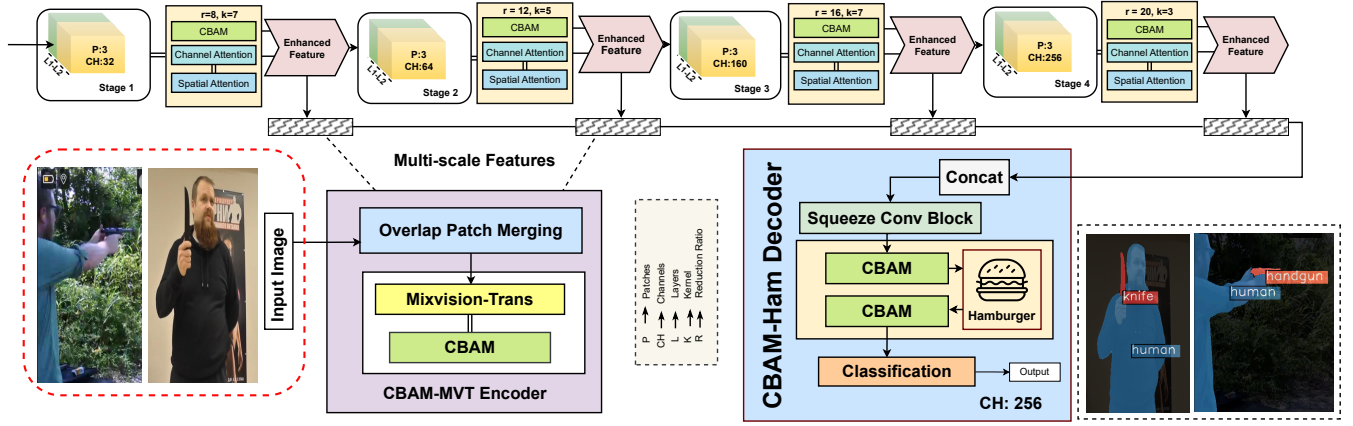


Fig. 2: ArmFormer Architecture Overview: The framework consists of a four-stage CBAM-enhanced MixVisionTransformer encoder backbone with progressive spatial resolution reduction and channel expansion, followed by a CBAM-integrated lightweight hamburger decoder head. The encoder processes input images through hierarchical stages (Stage 1-4) with progressive patch embeddings (P-3) and dual attention mechanisms (CBAM with Channel and Spatial Attention). Multi-scale features from all stages undergo overlap patch merging and are fused through the CBAM-MVT Encoder with Mixvision-Trans and CBAM modules. The decoder employs a squeeze convolution block, dual CBAM modules configured as a hamburger structure, and a classification head to generate pixel-level segmentation predictions for multi-class weapon detection.

efficient self-attention computation. The transformer blocks within each stage utilize efficient self-attention mechanisms that model long-range dependencies, combined with mix-FFN (Mix Feed-Forward Network) [58] layers that incorporate depth-wise convolutions for position encoding. Following each transformer stage, CBAM modules with uniform parameters (reduction ratio $r = 16$, spatial kernel size $k = 7$) are strategically inserted to enhance the extracted features through channel and spatial attention refinement.

The hierarchical design enables the encoder to capture both fine-grained local details from early stages and high-level semantic information from deeper stages, producing multi-scale feature representations $\{F_1, F_2, F_3, F_4\}$ that are essential for accurate dense prediction. The CBAM integration at each stage allows the network to focus on weapon-specific features while maintaining computational efficiency, with the attention mechanisms learning to emphasize relevant spatial regions and informative feature channels for multi-class weapon discrimination.

3) *CBAM-Integrated Decoder Head*: The decoder head employs a lightweight hamburger architecture with strategic dual CBAM integration to efficiently aggregate multi-scale features and generate pixel-level segmentation predictions. Multi-scale features $\{F_1, F_2, F_3, F_4\}$ from all encoder stages undergo overlap patch merging, where features are first aligned to unified spatial dimensions through bilinear interpolation, then concatenated along the channel dimension to form a comprehensive multi-scale representation.

The hamburger module performs efficient global context modeling through matrix decomposition, enabling the decoder to capture long-range dependencies with significantly reduced computational overhead compared to full self-attention mechanisms.

This is achieved by decomposing the feature map into lower-rank matrices that preserve global contextual information while maintaining linear complexity with respect to spatial resolution.

CBAM attention is applied at two strategic positions within the decoder pathway. The pre-hamburger CBAM (CBAM₁) refines the concatenated multi-scale features by enhancing channel-wise relationships and emphasizing informative spatial regions before global context aggregation. Following the hamburger module, post-hamburger CBAM (CBAM₂) further enhances the globally-contextualized features by applying sequential attention refinement. The complete decoder processing pipeline is formulated as:

$$S = \text{Conv}_{cls}(\text{CBAM}_2(\text{Ham}(\text{CBAM}_1(\text{Concat}(F_1, F_2, F_3, F_4)))))) \quad (4)$$

where Ham denotes the hamburger module for global context aggregation, Conv_{cls} is the classification convolution layer, and $S \in \mathbb{R}^{N \times H \times W}$ represents the final segmentation prediction with N classes including handgun, rifle, knife, revolver, and human.

The complete ArmFormer architecture integrates all components through end-to-end training using pixel-wise cross-entropy loss [59]:

$$\mathcal{L} = -\frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \sum_{n=1}^N y_{ij}^{(n)} \log(\hat{y}_{ij}^{(n)}) \quad (5)$$

where $y_{ij}^{(n)}$ represents the ground truth label and $\hat{y}_{ij}^{(n)}$ denotes the predicted probability for class n at spatial location (i, j) . This end-to-end optimization allows joint learning of attention mechanisms and segmentation performance, enabling

the network to adaptively focus on task-specific features tailored for multi-class weapon detection while maintaining the lightweight computational profile essential for real-time edge deployment.

V. EXPERIMENTAL RESULTS

This section focuses on conducting comprehensive experiments to evaluate the performance of our proposed ArmFormer against state-of-the-art semantic segmentation models across multiple dimensions including segmentation accuracy, computational efficiency, and per-class performance. The evaluation encompasses both lightweight models designed for edge deployment (CGNet [60], HrNet [61]) and heavyweight architectures optimized for accuracy (EncNet [62], ICNet [63], Uppernet_swin [64]), providing a thorough assessment of ArmFormer’s competitive positioning across the efficiency-accuracy spectrum.

A. Comparison with State-of-the-Art Methods

Table II presents the comprehensive performance comparison of our proposed ArmFormer against eight state-of-the-art semantic segmentation architectures. The evaluation metrics include mean Intersection over Union (mIoU), mean Accuracy (mAcc), mean F-score (mFscore), computational complexity (FLOPs in Giga operations), model parameters (in Millions), and inference speed (Frames Per Second).

TABLE II
PERFORMANCE COMPARISON

Model	FLOPs (G)	Params (M)	Speed (FPS)	mIoU (%)	mAcc (%)	mFscore (%)
ArmFormer (Ours)	4.886	3.66	82.26	80.64	88.28	89.13
CGNet [60]	3.452	0.493	90.49	64.30	74.96	77.08
ICNet [63]	15.434	47.52	140.88	74.16	82.41	84.63
Segmenter [46]	12.266	6.685	74.15	31.46	41.46	50.86
Uppernet_swin [64]	236.0	58.94	38.97	70.90	79.53	82.56
PspNet [65]	18.14	4.571	56.96	66.20	75.12	77.95
EncNet [62]	54.56	12.52	90.78	77.65	81.65	80.63
HrNet [61]	6.52	1.87	64.92	69.24	78.31	81.55

Table III provides detailed per-class IoU performance across all weapon categories and human detection, revealing the model’s capability to handle diverse object types in security-critical scenarios.

TABLE III
PER-CLASS IOU PERFORMANCE COMPARISON (%)

Model	Handgun	Human	Knife	Rifle	Revolver
ArmFormer (Ours)	82.24	67.22	80.81	83.87	80.43
CGNet [60]	71.13	32.92	66.42	72.67	61.90
ICNet [63]	78.72	50.12	76.83	77.46	75.65
Segmenter [66]	33.19	7.41	59.24	11.56	22.02
Uppernet_swin [64]	74.88	52.10	73.56	74.23	66.52
PspNet [65]	81.82	55.81	77.45	29.80	73.18
EncNet [62]	83.31	58.77	81.83	80.17	72.12
HrNet [61]	72.78	58.08	61.83	70.01	69.08

B. Performance Analysis and Model Superiority

ArmFormer achieves state-of-the-art performance with 80.64% mIoU and 89.13% mFscore, outperforming all baseline models including heavyweight architectures. Our model

surpasses EncNet [62] (77.65% mIoU) despite using 11.2× fewer FLOPs (4.886G vs 54.56G) and 3.4× fewer parameters (3.66M vs 12.52M), and significantly exceeds Uppernet_swin [64] (70.90% mIoU) while being 48× more computationally efficient. With 82.26 FPS inference speed, ArmFormer maintains real-time processing capabilities essential for security applications. While ICNet [63] achieves higher FPS (140.88), it sacrifices 6.48% mIoU, and CGNet [60] prioritizes speed (90.49 FPS) but underperforms by 16.34% mIoU, making both unsuitable for accurate threat assessment.

Table III demonstrates ArmFormer’s robust multi-class detection, achieving the highest performance in four out of five classes: Handgun (82.24%), Rifle (83.87%), Knife (80.81%), and Revolver (80.43%). Unlike EncNet [62] which shows inconsistent performance across categories, ArmFormer maintains consistently high accuracy across all weapon types, critical for comprehensive threat assessment. The comparison with Segmenter [66], another transformer-based approach, validates our CBAM integration strategy, with ArmFormer achieving 49.18% mIoU improvement due to domain-specific attention mechanisms.

The compact footprint of 4.886G FLOPs and 3.66M parameters positions ArmFormer as the optimal solution for edge deployment on resource-constrained devices including embedded GPUs, mobile processors, and AI accelerators in surveillance infrastructure. Unlike heavyweight models requiring server-grade hardware or ultra-lightweight models that sacrifice accuracy, ArmFormer achieves the critical balance necessary for practical real-world deployment in distributed security systems.

C. Qualitative Results Analysis

Figure 3 presents qualitative segmentation results that visually demonstrate the superior performance of our proposed ArmFormer compared to state-of-the-art baselines across diverse weapon detection scenarios. The visual comparison reveals critical insights into model capabilities for handling real-world security challenges.

In the revolver detection case (Figure 3 top row), ArmFormer produces the most accurate segmentation with complete object coverage and precise boundary delineation. While models like HrNet [61], ICNet [63], and PspNet [65] successfully detect the revolver, they exhibit incomplete segmentation with missing regions and fragmented predictions. Notably, Segmenter [66] completely fails to detect the revolver, producing an entirely black mask, which aligns with its extremely poor quantitative performance. Our ArmFormer achieves pixel-perfect segmentation with coherent object structure, demonstrating superior performance on the revolver class.

For rifle detection (Figure 3 middle row), the qualitative results definitively establish ArmFormer’s superiority. The rifle case is particularly challenging due to the elongated object structure and potential occlusions. While EncNet [62] and Uppernet_swin [64] show reasonable detection, they produce noisy predictions with boundary artifacts. Segmenter [66] again demonstrates catastrophic failure with minimal detection



Fig. 3: Qualitative comparison of segmentation results across different models on challenging weapon detection scenarios. The figure presents three representative test cases: revolver detection (top row), rifle segmentation (middle row), and handgun identification (bottom row). Each column represents predictions from different models (HrNet [61], ICNet [63], PspNet [65], Uppernet_swin [64], Segmenter [66], CGNet [60], EncNet [62], and our ArmFormer along with the original image). Ground truth annotations and background regions are shown with corresponding color coding for each weapon category. ArmFormer (rightmost column) consistently produces the most accurate pixel-level segmentation with precise object boundaries and minimal false predictions across all weapon categories.

capability. CGNet [60], despite its high inference speed, shows incomplete segmentation with missing rifle components. ArmFormer produces the most accurate and complete rifle segmentation with clean boundaries and comprehensive object coverage, directly correlating with its superior quantitative performance compared to competitors.

In the handgun detection case (Figure 3 bottom row), ArmFormer continues to demonstrate exceptional segmentation quality with precise boundary localization and complete object coverage. Most competing method is EncNet [62] which achieve reasonable handgun detection but suffer from boundary imprecision and occasional fragmentation. Uppernet_swin [64] shows competitive performance with relatively clean segmentation. However, Segmenter [66] maintains its poor performance pattern, failing to provide meaningful segmentation results. CGNet [60] produces acceptable results but with some boundary artifacts. ArmFormer’s handgun segmentation exhibits the highest fidelity with sharp, accurate boundaries and consistent object structure, reinforcing its robust performance across all weapon categories.

The qualitative analysis reveals a consistent pattern: ArmFormer maintains superior segmentation quality across all weapon categories and complexity levels. Unlike baseline models that show inconsistent performance across different scenarios, our proposed architecture demonstrates robust generalization enabled by strategic CBAM integration and transformer-based global context modeling. The visual results validate that ArmFormer’s architectural design effectively addresses the fundamental challenges in weapon segmentation, including multi-scale feature representation, precise boundary localization, and reliable multi-class discrimination essential for security-critical applications.

VI. ABLATION STUDY

We conducted three comprehensive ablation studies on our proposed ArmFormer model to validate the effectiveness of our design choices and demonstrate the superiority of our approach. The ablation studies systematically evaluate: (1) the impact of FPN neck integration for multi-scale feature fusion, (2) the effect of lightweight CBAM configurations for computational efficiency, and (3) the influence of adaptive stage-specific CBAM parameters.

TABLE IV
ABLATION STUDY RESULTS: ARMFORMER VS. ALTERNATIVE CONFIGURATIONS

Model Configuration	mIoU (%)	mFscore (%)	FPS
ArmFormer	80.64	89.13	82.26
w/ FPN Neck	79.92	88.70	78.85
w/ Lightweight CBAM	79.74	88.57	84.28
w/ Adaptive CBAM	78.05	87.54	81.15

TABLE V
PER-CLASS IOU PERFORMANCE (%)

Configuration	Bg	Gun	Human	Knife	Rifle	Rev
ArmFormer	89.29	82.24	67.22	80.81	83.87	80.43
w/ FPN Neck	89.01	79.65	67.97	81.60	82.38	78.94
w/ Lightweight	89.46	79.88	66.90	82.87	81.32	78.04
w/ Adaptive	87.20	79.30	66.77	79.95	76.52	78.53

Our comprehensive ablation studies demonstrate the superiority of our proposed ArmFormer design and validate our architectural choices. The results clearly show that our base ArmFormer achieves the best overall performance with

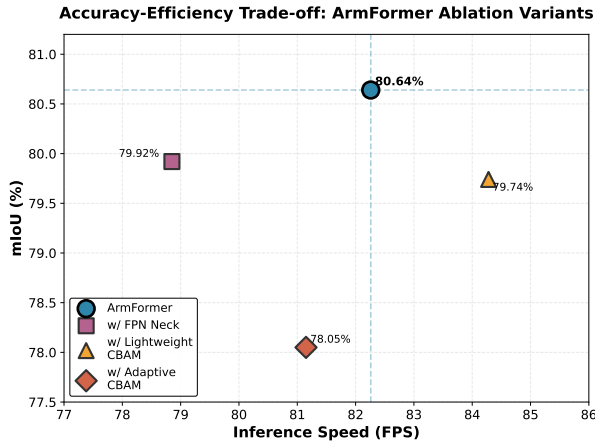


Fig. 4: Accuracy-efficiency trade-off visualization for ArmFormer ablation variants. The scatter plot illustrates the relationship between inference speed (FPS) and segmentation accuracy (mIoU) across different configurations. Our base ArmFormer (blue circle) achieves the optimal balance, positioned in the upper-right region with highest mIoU (80.64%) and competitive FPS (82.26). The dashed reference lines intersect at ArmFormer’s performance point, demonstrating its superior positioning in the accuracy-efficiency space.

80.64% mIoU and 89.13% mFscore, outperforming all ablation variants. Figure 4 visualizes the accuracy-efficiency trade-off across all configurations, revealing critical insights into the design space exploration.

FPN Neck Integration Analysis. As illustrated in Figure 4, the FPN neck variant (red square) falls into the lower-left quadrant, exhibiting both reduced accuracy and slower inference. Adding FPN neck integration results in a performance decline from 80.64% to 79.92% mIoU (-0.72%) and reduced inference speed (82.26 to 78.85 FPS). While FPN theoretically provides better multi-scale feature fusion, our results demonstrate that the additional computational overhead outweighs the benefits, validating our decision to exclude FPN from the main ArmFormer architecture. The graph clearly shows that this configuration moves away from the optimal trade-off point.

Lightweight CBAM Configuration. The lightweight variant (orange triangle) appears in the upper-right region of Figure 4, achieving the highest FPS (84.28) among all configurations. This variant with higher reduction ratios (32) and smaller kernels (3) achieves 79.74% mIoU while improving computational efficiency. Although this configuration offers better speed performance, it sacrifices 0.90% mIoU compared to our main model. The visualization reveals that while lightweight CBAM provides marginal speed gains (+2.02 FPS), the accuracy degradation makes it suboptimal for security-critical applications where precision is paramount. Our balanced CBAM parameters (reduction ratio: 16, kernel size: 7) provide the optimal accuracy-efficiency trade-off, as evidenced by ArmFormer’s superior positioning in the accuracy-efficiency

space.

Adaptive CBAM Parameters. The adaptive configuration (brown diamond) demonstrates the worst performance in Figure 4, positioned in the lower region with the lowest mIoU (78.05%, -2.59% degradation). This counterintuitive result indicates that overly complex adaptive mechanisms can hinder optimization and lead to suboptimal performance, particularly evident in rifle detection (83.87% \rightarrow 76.52%). The graph visualization definitively shows that adaptive stage-specific parameters do not improve the accuracy-efficiency trade-off, validating our design choice of uniform CBAM parameters across stages. The performance point’s position away from the optimal frontier confirms that architectural simplicity combined with strategic attention placement yields superior results.

Our ablation studies conclusively prove that the proposed ArmFormer with uniform CBAM integration achieves superior performance compared to all alternative configurations. The consistent attention mechanisms across stages, optimal CBAM parameters, and streamlined architecture without FPN neck represent the most effective design for weapon detection segmentation.

VII. CONCLUSION

This paper presented ArmFormer, a lightweight transformer-based semantic segmentation framework for real-time multi-class weapon detection in edge environments. By strategically integrating CBAM with MixVisionTransformer architecture, ArmFormer achieves state-of-the-art performance with 80.64% mIoU and 89.13% mFscore at 82.26 FPS, outperforming heavyweight models requiring up to 48 \times more computation. With only 4.886G FLOPs and 3.66M parameters, ArmFormer is optimally suited for deployment on portable security cameras, programmable surveillance systems, autonomous drones, and embedded AI accelerators, enabling real-time threat assessment without cloud dependency. Comprehensive ablation studies validate that uniform CBAM integration across encoder stages yields superior accuracy-efficiency trade-offs. Despite these achievements, limitations include performance variability under extreme lighting conditions, challenges with heavily occluded weapons, and the need for larger-scale diverse training datasets. Future work will explore model quantization and pruning for ultra-low-power scenarios, multi-modal fusion with thermal and depth sensors for enhanced detection under challenging conditions, real-time video stream processing with temporal consistency, and federated learning strategies for privacy-preserving collaborative training across distributed edge devices.

REFERENCES

- [1] U.S. Centers for Disease Control and Prevention, “Injury and violence data, 2024-2025,” Online, CDC, 2024, u.S. Centers for Disease Control and Prevention.
- [2] A. Glowacz, M. Kmiec, and A. Dziech, “Visual detection of knives in security applications using active appearance models,” *Multimedia Tools and Applications*, vol. 74, pp. 4253–4267, 2015. [Online]. Available: <https://doi.org/10.1007/s11042-013-1537-2>

- [3] X. Dong *et al.*, “Real-time weapon detection using yolov8 for enhanced safety,” *arXiv preprint arXiv:2410.19862*, 2024.
- [4] DataCamp, “Yolo object detection explained: A beginner’s guide,” Online, 2024, available at DataCamp. [Online]. Available: <https://www.datacamp.com>
- [5] R. K. Tiwari and G. K. Verma, “A computer vision based framework for visual gun detection using harris interest point detector,” *Procedia Computer Science*, vol. 54, pp. 703–712, 2015.
- [6] A. S. V. Rao *et al.*, “An efficient weapon detection system using nsgcudnn classifier in surveillance,” *Expert Systems with Applications*, vol. 255, p. 124800, 2024.
- [7] M. S. Nadeem *et al.*, “A comprehensive study towards high-level approaches for weapon detection using classical machine learning and deep learning methods,” *Expert Systems with Applications*, vol. 212, p. 118746, 2023.
- [8] R. Debnath and M. K. Bhowmik, “A comprehensive survey on computer vision-based concepts, methodologies, analysis and applications for automatic gun/knife detection,” *Journal of Visual Communication and Image Representation*, vol. 78, p. 103165, 2021.
- [9] K. Akhila and K. R. Ahmed, “Real time deep learning weapon detection techniques for mitigating lone wolf attacks,” *arXiv preprint arXiv:2405.14148*, 2024.
- [10] P. Y. Ingle and Y.-G. Kim, “Real-time abnormal object detection for video surveillance in smart cities,” *Sensors*, vol. 22, no. 10, p. 3862, 2022.
- [11] IBM Research, “What is semantic segmentation?” IBM Think Topics, 2024, available online. [Online]. Available: <https://www.ibm.com>
- [12] Keymakr, “Semantic segmentation vs object detection: Understanding the differences,” Online, 2024, available at Keymakr. [Online]. Available: <https://www.keymakr.com>
- [13] E. Xie *et al.*, “Segformer: Simple and efficient design for semantic segmentation with transformers,” in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 12077–12090, neurIPS 2021.
- [14] Hugging Face, “Segformer: Simple and efficient design for semantic segmentation with transformers,” Documentation, 2024, available at Hugging Face. [Online]. Available: <https://huggingface.co>
- [15] A. Egiazarov, V. Mavroedis, F. M. Zennaro, and K. Vishi, “Firearm detection and segmentation using an ensemble of semantic neural networks,” *arXiv preprint arXiv:2003.00805*, 2020.
- [16] V. Kaya, S. Tuncer, and A. Baran, “Detection and classification of different weapon types using deep learning,” *Applied Sciences*, vol. 11, no. 16, p. 7535, 2021.
- [17] T. Santos, H. Oliveira, and A. Cunha, “Systematic review on weapon detection in surveillance footage through deep learning,” *Computer science review*, vol. 51, p. 100612, 2024.
- [18] H. Bai, H. Mao, and D. Nair, “Dynamically pruning segformer for efficient semantic segmentation,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, May 2022, pp. 3298–3302.
- [19] M. Contributors, “MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark,” <https://github.com/open-mmlab/mms Segmentation>, 2020.
- [20] Z. Liu *et al.*, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021, pp. 10012–10022.
- [21] A. Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [22] S. Woo *et al.*, “Cbam: Convolutional block attention module,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2018, pp. 3–19.
- [23] Internet Movie Firearms Database, “Category: Gun,” Online, 2024, Internet Movie Firearms Database (IMFDB). [Online]. Available: <https://www.imfdb.org/wiki/Category:Gun>
- [24] Y. Li *et al.*, “A microwave measurement system for metallic object detection using swept-frequency radar,” in *Millimetre Wave and Terahertz Sensors and Technology*, vol. 7117. International Society for Optics and Photonics, 2008, p. 71170K.
- [25] D. M. Sheen, D. L. McMakin, and T. E. Hall, “Three-dimensional millimeter-wave imaging for concealed weapon detection,” *IEEE Transactions on Microwave Theory and Techniques*, vol. 49, no. 9, pp. 1581–1592, 2001.
- [26] Z. Xue and R. S. Blum, “Concealed weapon detection using color image fusion,” in *Proceedings of the 6th International Conference on Information Fusion*, vol. 1. IEEE, 2003, pp. 622–627.
- [27] Transportation Security Administration, “Tsa intercepts 6,678 firearms at airport security checkpoints in 2024,” Press Release, January 2025, Washington, DC: Transportation Security Administration. [Online]. Available: <https://www.tsa.gov/news/press/releases/2025/01/15/tsa-intercepts-6678-firearms-airport-security-checkpoints-2024>
- [28] R. Gesick, C. Saritac, and C.-C. Hung, “Automatic image analysis process for the detection of concealed weapons,” in *Proceedings of the 5th Annual Workshop on Cyber Security and Information Intelligence Research: Cyber Security and Information Intelligence Challenges and Strategies*, 2009, pp. 1–4.
- [29] G. Flitton, T. P. Breckon, and N. Megherbi, “A comparison of 3d interest point descriptors with application to airport baggage object detection in complex ct imagery,” *Pattern Recognition*, vol. 46, no. 9, pp. 2420–2436, 2013.
- [30] X. Zelong *et al.*, “Automatic detection of concealed pistols using passive millimeter wave imaging,” in *2015 IEEE International Conference on Imaging Systems and Techniques (IST)*. IEEE, 2015, pp. 1–4.
- [31] A. Kambhatla and K. R. Ahmed, “Firearm detection using deep learning,” 2023.
- [32] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection,” vol. 1, pp. 886–893, 2005.
- [33] I. Darker *et al.*, “Can cctv reliably detect gun crime?” in *2007 41st Annual IEEE International Carnahan Conference on Security Technology*. IEEE, 2007, pp. 264–271.
- [34] L. Wang, S. Guo, W. Huang, and Y. Qiao, “Places205-vggnet models for scene recognition,” *arXiv preprint arXiv:1508.01667*, 2015.
- [35] S. Targ, D. Almeida, and K. Lyman, “Resnet in resnet: Generalizing residual architectures,” *arXiv preprint arXiv:1603.08029*, 2016.
- [36] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, others, and L. Farhan, “Review of deep learning: concepts, cnn architectures, challenges, applications, future directions,” *Journal of Big Data*, vol. 8, no. 1, p. 53, 2021.
- [37] P. Jiang, D. Ergu, F. Liu, Y. Cai, and B. Ma, “A review of yolo algorithm developments,” *Procedia Computer Science*, vol. 199, pp. 1066–1073, 2022.
- [38] J. Redmon and A. Farhadi, “Yolov3: An incremental improvement,” *arXiv preprint arXiv:1804.02767*, 2018.
- [39] A. Bochkovskiy, C. Y. Wang, and H. Y. M. Liao, “Yolov4: Optimal speed and accuracy of object detection,” *arXiv preprint arXiv:2004.10934*, 2020.
- [40] R. Debnath and M. K. Bhowmik, “A comprehensive survey on computer vision-based weapon detection,” *Journal of Visual Communication*, 2021.
- [41] A. Krizhevsky *et al.*, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, vol. 25. Curran Associates, Inc., 2012, pp. 1097–1105.
- [42] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [43] N. Carion *et al.*, “End-to-end object detection with transformers,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 213–229.
- [44] J. Long *et al.*, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015, pp. 3431–3440.
- [45] O. Ronneberger *et al.*, “U-net: Convolutional networks for biomedical image segmentation,” in *Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [46] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, “Segmenter: Transformer for semantic segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 7262–7272.
- [47] W. Wang *et al.*, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021, pp. 568–578.
- [48] K. Han *et al.*, “A survey on vision transformer,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [49] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2018, pp. 7132–7141.

- [50] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019, pp. 510–519.
- [51] F. Wang *et al.*, "Residual attention network for image classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017.
- [52] Y. Chen *et al.*, "Attention mechanisms in computer vision: A survey," *Computational Visual Media*, 2022.
- [53] M. Z. Alom *et al.*, "A state-of-the-art survey on deep learning theory and architectures," *Electronics*, 2019.
- [54] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.
- [55] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, and C. Feichtenhofer, "Sam 2: Segment anything in images and videos," *arXiv preprint arXiv:2408.00714*, 2024. [Online]. Available: <https://arxiv.org/abs/2408.00714>
- [56] T. Islam, T. T. Sarker, K. R. Ahmed, C. B. Rankrape, and K. Gage, "Weedswin hierarchical vision transformer with sam-2 for multi-stage weed detection and classification," *Scientific Reports*, vol. 15, no. 1, p. 23274, 2025.
- [57] X. Yu, J. Wang, Y. Zhao, and Y. Gao, "Mix-vit: Mixing attentive vision transformer for ultra-fine-grained visual categorization," *Pattern Recognition*, vol. 135, p. 109131, 2023.
- [58] N. Aftabi, N. Moradi, and F. Mahroo, "Feed-forward neural networks as a mixed-integer program," *Engineering with Computers*, pp. 1–19, 2025.
- [59] A. Mao, M. Mohri, and Y. Zhong, "Cross-entropy loss functions: Theoretical analysis and applications," in *International conference on Machine learning*. pmlr, 2023, pp. 23 803–23 828.
- [60] T. Wu, S. Tang, R. Zhang, J. Cao, and Y. Zhang, "Cgnet: A light-weight context guided network for semantic segmentation," *IEEE Transactions on Image Processing*, vol. 30, pp. 1169–1179, 2020.
- [61] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *CVPR*, 2019.
- [62] H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal, "Context encoding for semantic segmentation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [63] H. Zhao, X. Qi, X. Shen, J. Shi, and J. Jia, "Icnet for real-time semantic segmentation on high-resolution images," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 405–420.
- [64] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, "Unified perceptual parsing for scene understanding," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 418–434.
- [65] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.
- [66] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, "Segmenter: Transformer for semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 7262–7272.