

# Track Any Anomalous Object: A Granular Video Anomaly Detection Pipeline

## Supplementary Material

### 679 6. Supplementary

680 We propose a straightforward and effective pipeline for fine-  
681 grained video anomaly detection. The core approach consists  
682 of two key steps: **1)** generating anomaly prompts based  
683 on object detection results from video sequences and refining  
684 these prompts using a robust filtering algorithm; **2)**  
685 applying a prompt-based segmentation model to produce  
686 accurate pixel-level anomaly masks. This pipeline enables  
687 efficient identification and segmentation of anomalous objects,  
688 significantly improving detection precision and temporal  
689 consistency. In the supplementary materials, we provide  
690 additional details on the following aspects:

- 691 • We offer a more comprehensive explanation of the im-  
692 plementation process for anomalous boxes extraction and  
693 SAM2 segmentation inference in Sec. 6.1.
- 694 • We provide a detailed explanation of the object-level  
695 evaluation metrics utilized in our experiments, empha-  
696 sizing their importance in assessing spatial and temporal  
697 anomaly detection performance in Sec. 6.2.
- 698 • We analyze the limitations of our proposed model and  
699 existing baselines, highlighting potential avenues for im-  
700 provement in Sec. 6.3.
- 701 • we present additional experimental visualization results,  
702 highlighting the instance segmentation performance of  
703 our model, including examples from the ShanghaiTech  
704 Campus dataset in Sec. 6.4.

### 705 6.1. Comprehensive Implementation Details

706 **Details of Anomalous Boxes Extraction.** Our object-level  
707 VAD algorithm [27] utilizes features such as speed, pose,  
708 and depth to detect anomalies. During the training phase,  
709 a probabilistic density model, such as k-nearest neighbors  
710 or Mahalanobis distance, is constructed based on normal  
711 behavioral attributes. In the testing phase, the probability  
712 density of each object’s features is calculated, where lower  
713 density values indicate greater deviation from normal be-  
714 havior, resulting in higher anomaly scores. For the test data,  
715 only speed and depth features are used to compute anomaly  
716 scores. These scores are then standardized using the corre-  
717 sponding speed and depth anomaly scores from the training  
718 data, enhancing the prominence of anomalous objects and  
719 making them easier to detect. The standardized speed and  
720 depth anomaly scores are subsequently summed to produce  
721 a final overall anomaly score. A threshold is applied to this  
722 score to effectively filter and identify anomalous objects and  
723 their corresponding bounding boxes.

724 To address overlapping anomaly boxes that may corre-  
725 spond to the same anomalous object, we calculate the Inter-

section over Union (IoU) for each pair of filtered boxes  $B_i$  726  
and  $B_j$  within the same frame. The IoU is computed as the 727  
ratio of the intersection area to the union area between two 728  
boxes, defined as: 729

$$\text{IoU}(B_i, B_j) = \frac{\text{Area}(B_i \cap B_j)}{\text{Area}(B_i \cup B_j)}. \quad 730$$

If the IoU exceeds a threshold of  $\tau = 0.3$ , the two boxes are 731  
deemed to represent the same anomalous object. In such 732  
cases, a new box  $B_{\text{new}}$  is created by merging the two, with 733  
its top-left corner coordinates given by  $\min(x_i^{\min}, x_j^{\min})$  and 734  
 $\min(y_i^{\min}, y_j^{\min})$ , and its bottom-right corner coordinates de- 735  
termined as  $\max(x_i^{\max}, x_j^{\max})$  and  $\max(y_i^{\max}, y_j^{\max})$ . This 736  
merging process consolidates overlapping boxes into a single 737  
bounding box that accurately represents the anomalous 738  
object. By iteratively applying this procedure across all 739  
frames, the algorithm ensures a non-redundant and consis- 740  
tent representation of anomalies, improving the overall ac- 741  
curacy and reliability of localization. 742

**Details of Segmentation Model Inference.** We utilize 743  
SAM2 as the prompt-based segmentation model to per- 744  
form instance segmentation for distinct anomalous objects 745  
in video clips. This process involves generating prompts 746  
for each object using a robust bounding box filtering al- 747  
gorithm, applied at fixed frame intervals. Each prompt is 748  
stored as a tuple  $\mathcal{T}_{\text{box}} = (f_i, b_j, \mathcal{L}_j)$ , where  $f_i$  represents the 749  
frame,  $b_j$  is the bounding box, and  $\mathcal{L}_j$  is the correspond- 750  
ing object label. These tuples ensure consistent and accu- 751  
rate tracking of anomalous objects across frames, maintain- 752  
ing both spatial and temporal coherence. By consolidat- 753  
ing bounding boxes and their associated labels into struc- 754  
tured prompts, the model effectively localizes and segments 755  
anomalies within dynamic video contexts. Instance seg- 756  
mentation in SAM2 is performed by providing the corre- 757  
sponding object labels and bounding box prompts as inputs. 758  
The resulting segmentation outputs are processed to serve 759  
different evaluation purposes. For pixel-level metrics, the 760  
instance segmentation results are transformed into binarized 761  
segmentation masks that highlight anomalous regions at the 762  
pixel level. For object-level metrics, each instance segmen- 763  
tation result is converted into a bounding box that encapsu- 764  
lates the segmented object. This dual-processing approach 765  
allows for comprehensive evaluation of both fine-grained 766  
anomaly localization and high-level object tracking, ensur- 767  
ing the robustness and accuracy of the proposed method. 768

### 769 6.2. Object-Level Evaluation Metrics

**RBDC.** The Region-Based Detection Criterion (RBDC) 770  
evaluates the spatial accuracy of anomaly detection by 771

772 quantifying the proportion of correctly matched predicted  
773 regions relative to the ground truth regions. For a predicted  
774 bounding box  $B_{\text{pred}}$  and a ground truth box  $B_{\text{gt}}$ , the Inter-  
775 section over Union (IoU) is computed as:

$$776 \quad \text{IoU} = \frac{\text{Area}(B_{\text{pred}} \cap B_{\text{gt}})}{\text{Area}(B_{\text{pred}} \cup B_{\text{gt}})}.$$

777 A match is deemed correct if  $\text{IoU} > \alpha$ , where  $\alpha$  is a prede-  
778 fined threshold (e.g.,  $\alpha = 0.1$ ). The RBDC score is defined  
779 as:

$$780 \quad \text{RBDC} = \frac{\text{Number of Correctly Matched Regions}}{\text{Total Number of Ground Truth Regions}}.$$

781 RBDC is essential for assessing spatial precision, ensuring  
782 that detected anomalies accurately align with ground truth  
783 regions, particularly in scenarios with complex or overlap-  
784 ping anomalies.

785 **TBDC.** The Track-Based Detection Criterion (TBDC) eval-  
786 uates the temporal consistency of anomaly detection by  
787 measuring the proportion of correctly tracked anomaly tra-  
788 jectories relative to the total number of ground truth tracks.  
789 A trajectory is considered correctly tracked if the IoU be-  
790 tween the predicted bounding box and the ground truth box  
791 exceeds the threshold  $\alpha$  in each frame of the track. The  
792 TBDC score is calculated as:

$$793 \quad \text{TBDC} = \frac{\text{Number of Correctly Tracked Anomaly Tracks}}{\text{Total Number of Ground Truth Tracks}}.$$

794 TBDC is critical for evaluating temporal robustness, cap-  
795 turing the model’s ability to maintain consistent anomaly  
796 detection across consecutive frames, which is particularly  
797 important in dynamic video scenarios.

798 Together, RBDC and TBDC provide a comprehensive  
799 framework for evaluating object-level anomaly detection,  
800 addressing both spatial precision and temporal coherence.  
801 These metrics are particularly well-suited for real-world ap-  
802 plications such as surveillance and autonomous systems,  
803 where accurate spatial localization and robust temporal  
804 tracking are paramount.

### 805 6.3. Limitations and Future Directions

806 **Limitations.** The performance of our model is closely tied  
807 to the robustness of the prompts provided to the prompt-  
808 based segmentation model, such as SAM2. These prompts  
809 heavily rely on the effectiveness of the object-level anomaly  
810 detection algorithm in assigning accurate anomaly scores to  
811 objects within anomalous frames. A precise anomaly detec-  
812 tion algorithm that effectively distinguishes between normal  
813 and anomalous objects generates higher-quality prompts,  
814 resulting in improved segmentation accuracy. However,  
815 false-positive prompts pose significant challenges. First,  
816 they can introduce cumulative tracking errors in SAM2,

leading to catastrophic forgetting of actual anomalous ob- 817  
jects. As these errors propagate across frames, the model 818  
may progressively lose its ability to detect critical anoma- 819  
lies, severely undermining its reliability. Second, false- 820  
positive prompts increase the computational burden dur- 821  
ing SAM2 inference. By prompting the model to pro- 822  
cess non-anomalous objects, they degrade inference effi- 823  
ciency, resulting in slower processing times and unneces- 824  
sary computational overhead. Addressing these issues is 825  
essential to ensure both the accuracy and efficiency of the 826  
proposed framework. This underscores the need to enhance 827  
the anomaly detection algorithm’s precision and robustness, 828  
thereby minimizing the impact of false-positive prompts 829  
and maximizing the overall performance of the system. 830

**Future Directions.** Existing video anomaly detection 831  
datasets primarily focus on frame-level and object-level 832  
anomalies, with pixel-level annotations being extremely 833  
limited. Among the few datasets that provide pixel-level 834  
annotations, these are often coarse, offering only rough out- 835  
lines of anomalous objects rather than precise contours. Ad- 836  
ditionally, current pixel-level annotations are predominantly 837  
binary masks, which pose significant challenges in scenar- 838  
ios where anomalous objects overlap, as binary masks fail to 839  
differentiate between overlapping objects, making accurate 840  
evaluation difficult. To address these limitations, we pro- 841  
pose adopting instance-level pixel annotations for anoma- 842  
lous objects. Instance-level annotations would uniquely 843  
identify each anomalous object at the pixel level, even in 844  
complex scenarios involving overlapping objects. This ap- 845  
proach would enhance the precision of pixel-level anomaly 846  
detection and enable more granular evaluations in video 847  
anomaly detection tasks. Furthermore, the adoption of 848  
instance-level pixel annotations would support the develop- 849  
ment of more robust algorithms capable of handling real- 850  
world scenarios, where anomalies often appear in complex 851  
and overlapping forms. By bridging the gap in current 852  
datasets, instance-level annotations could serve as a criti- 853  
cal foundation for advancing video anomaly detection and 854  
promoting consistent benchmarking across future methods. 855

### 856 6.4. More Visualization about Experiments

857 We provide additional visualization results from compara- 858  
tive experiments, showcasing detailed visual comparisons 859  
between our method and four baselines—SimpleNet [17], 860  
DRAEM [36], DDAD [21], and AnomalyCLIP [38]—on 861  
three selected video clips from the UCSD Ped2 dataset (see 862  
Fig. 7–9). Binary segmentation masks are employed for 863  
consistency and clarity, and the results clearly demonstrate 864  
that our model significantly outperforms the baseline meth- 865  
ods. Additionally, we present instance segmentation results 866  
of our model on the ShanghaiTech Campus dataset, further 867  
illustrating its effectiveness and robustness in segmenting 868  
anomalous objects at the instance level (see Fig. 10).

Figure 7. Visual comparisons on the UCSD Ped2 dataset. Red masks represent the binary segmentation results.

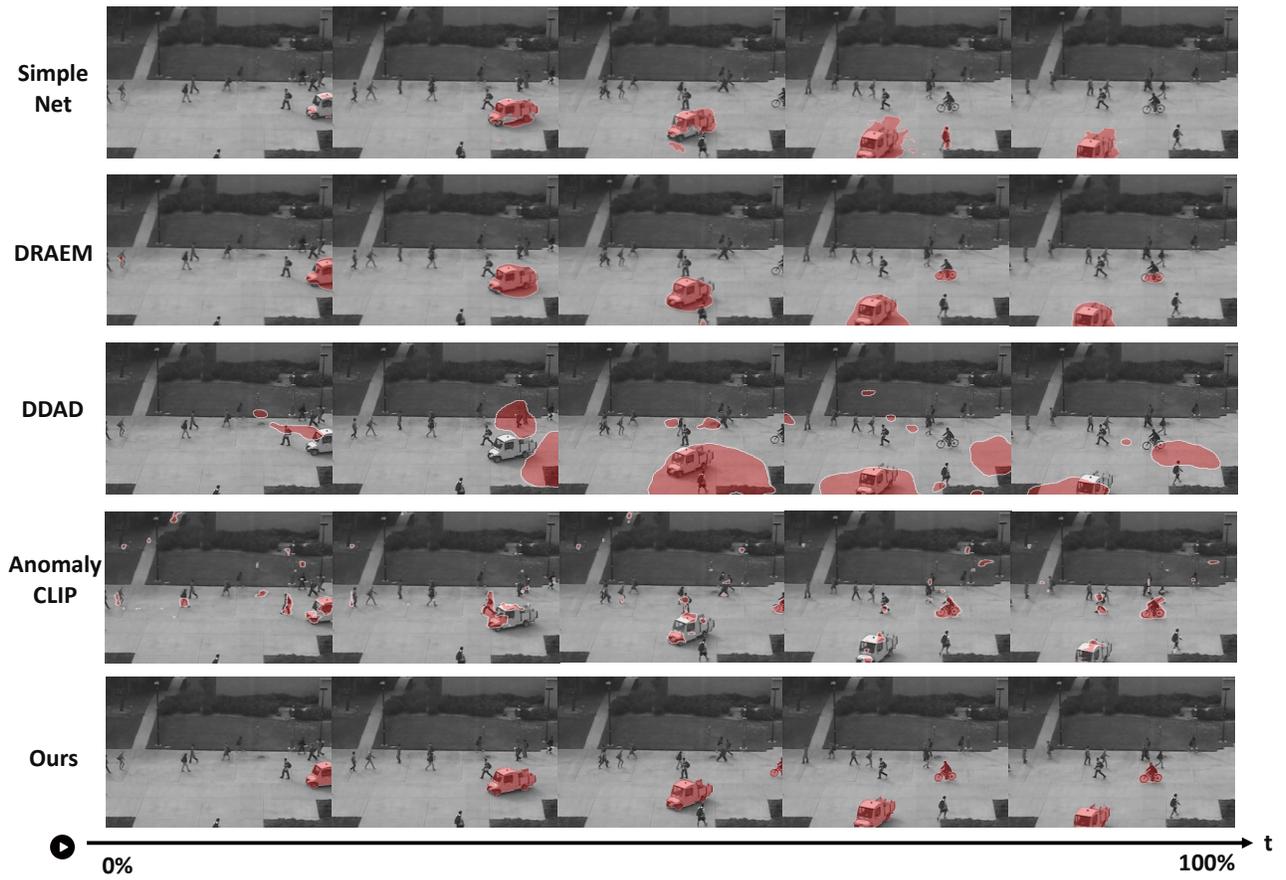


Figure 8. Visual comparisons on the UCSD Ped2 dataset. Red masks represent the binary segmentation results.

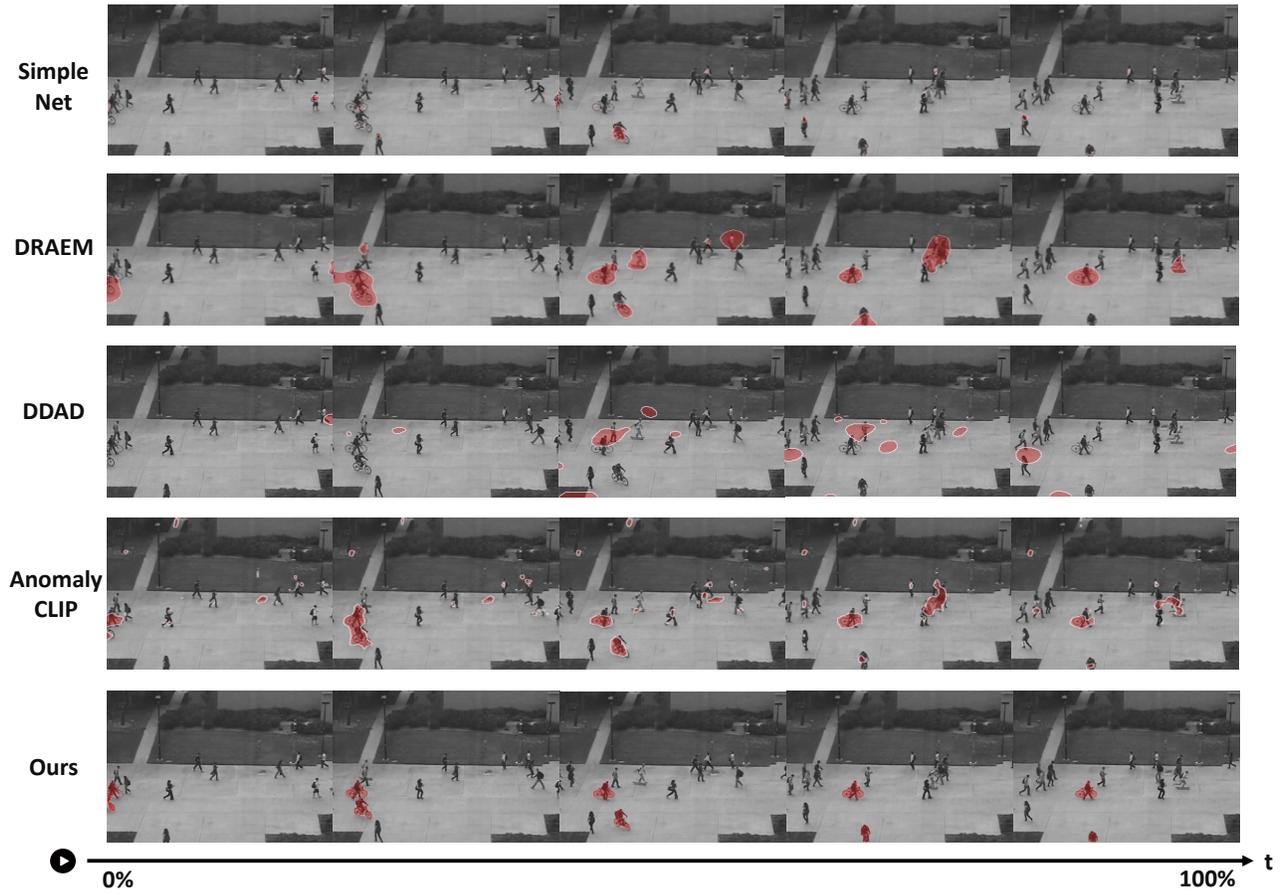


Figure 9. Visual comparisons on the UCSD Ped2 dataset. Red masks represent the binary segmentation results.

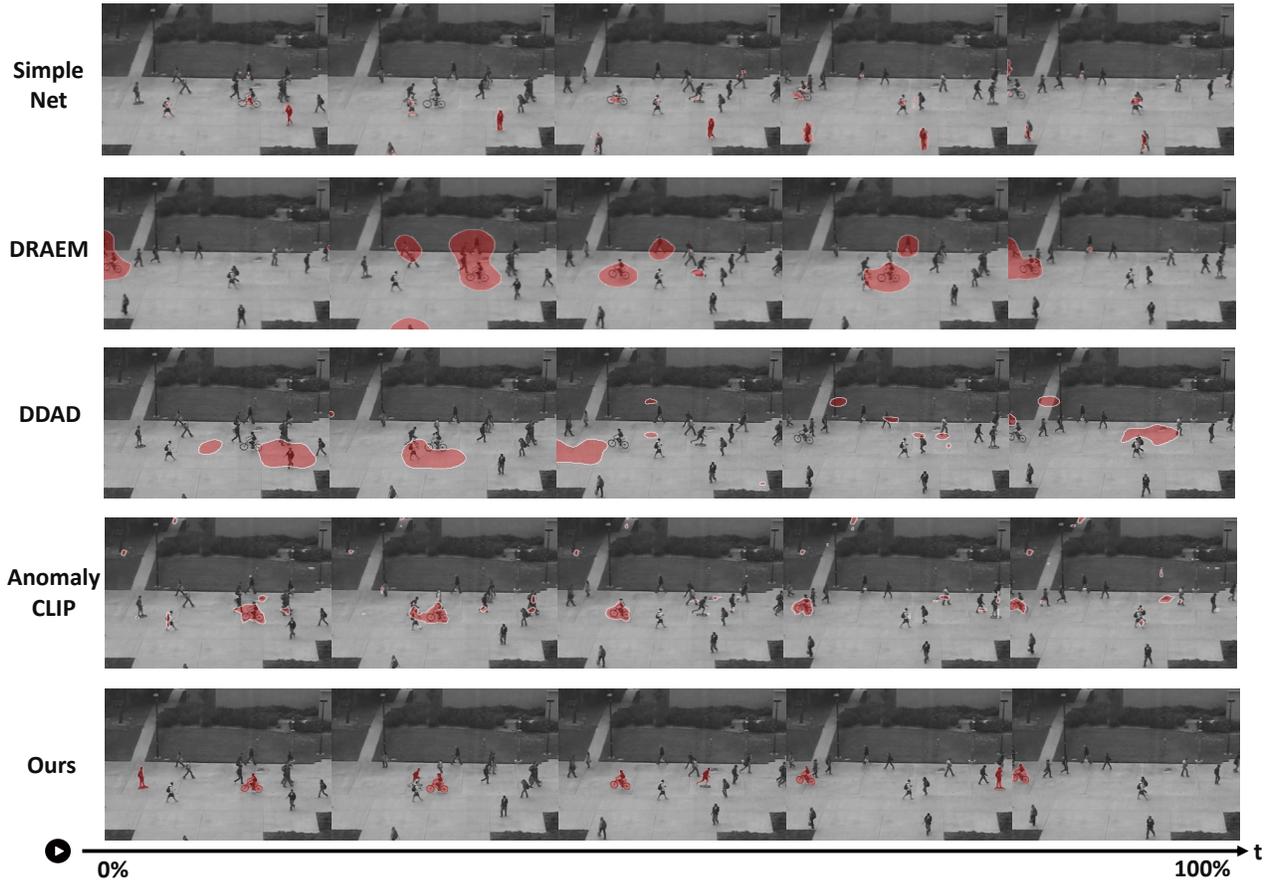


Figure 10. Instance segmentation visualization results on the ShanghaiTech Campus dataset. Masks in different colors represent the segmentation results for distinct instances.

