

ProstaTD: Bridging Surgical Triplet from Classification to Fully Supervised Detection

Paper ID: 1839

Presenter: Yiliang Chen

Background



- **Definition:** A surgical triplet is made up of three components which is the **instrument**, the **action**, and the **target**, which is defined by medical experts.
- **Objective:** The goal is to automatically **recognize** these triplets [tool, action, target] from each frame of surgical videos.

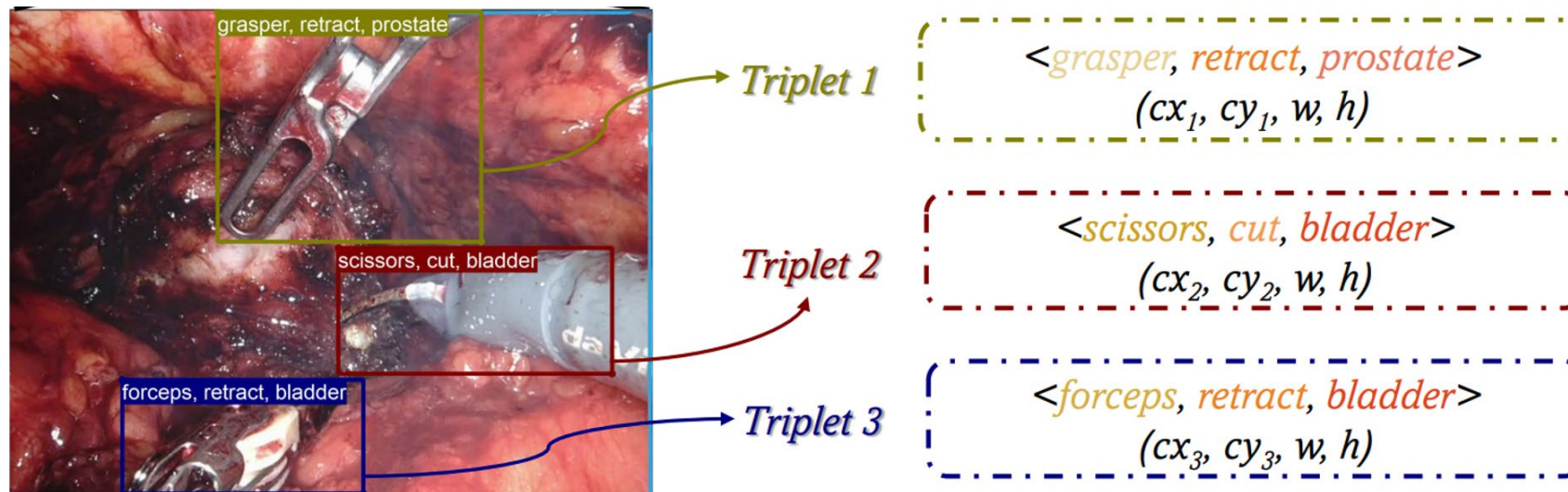




Table 1: Comparison of existing surgical datasets with ProstaTD, highlighting the unique attributes required for surgical triplet detection. LC and RARP denote Laparoscopic Cholecystectomy and Robot-Assisted Radical Prostatectomy. Full BBox indicates that every frame is fully annotated with bounding boxes. Complete Surgery denotes that the dataset as a whole provides coverage for the entire surgical procedure. Triplet Boundary refers to the temporal start and end range of a complete surgical triplet. Multiple Sources indicates that the data is collected from various hospital.

Dataset	Task	Attributes						Statistics	
		Supervised Detection	Full BBox	Triplet-like Structure	Triplet Boundary	Multiple Sources	Complete Surgery	No. Instances	No. Triplets
Cholec80-locations	LC	✓	✓				✓	6,471	–
CholecTrack20	LC	✓	✓				✓	65,200	–
ESAD	RARP	✓	✓				✓	46,753	–
PSI-AVA	RARP	✓					✓	5,804	–
CholecQ	LC	✓	✓	✓				14,480	17
CholecT45	LC			✓			✓	146,394	100
CholecT50	LC			✓			✓	161,988	100
ProstaTD (Ours)	RARP	✓	✓	✓	✓	✓	✓	196,490	89



In summary, our contributions are fourfold:

- We introduce a new task with a new dataset for fully supervised surgical triplet detection at the procedure level. To the best of our knowledge, our ProstaTD is the *largest* surgical dataset with instance-level annotations.
- We construct a multi-institutional surgical triplet dataset featuring detailed labels (precise bounding boxes and standardized triplet boundaries) in more complex surgical scenarios.
- We release two open-source annotation tools, which are the first specifically tailored for surgical triplet annotation, along with an open-source evaluation toolkit for benchmarking surgical triplet detection, providing a foundation for surgical triplet analysis across diverse surgical procedures.
- We introduce the first benchmark for fully supervised surgical triplet detection, providing our tailored method as a baseline for comparison.

Dataset Structures

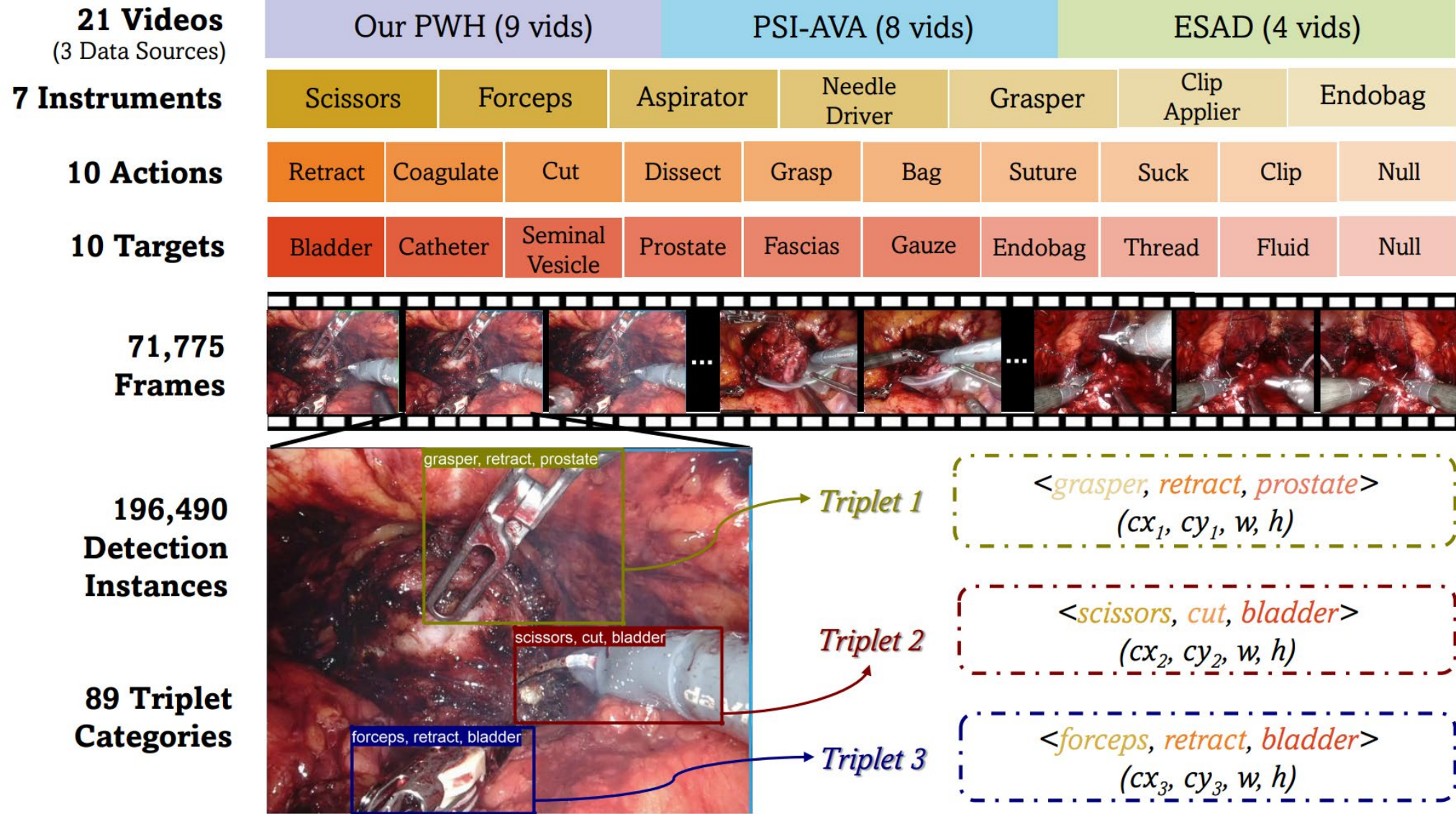


Table 4: Detection performance on the ProstaTD dataset for I, V, T, and IVT components. We report mAP at IoU thresholds of 50% (“50”) and 50:95 (“95”), together with inference speed (FPS). All results are reported as mean $_{\pm\text{std}}$ (%) over 5-fold cross-validation. Experiments are conducted with input size 640×640 on a single NVIDIA RTX 4090 GPU. Bold values with light green background indicate the best results, and underlined values with light purple background indicate the second-best.

Method	$mAP_I(\%) \uparrow$		$mAP_V(\%) \uparrow$		$mAP_T(\%) \uparrow$		$mAP_{IVT}(\%) \uparrow$		FPS \uparrow
	50	95	50	95	50	95	50	95	
Tripnet-Det*	1.6 _{0.4}	–	0.6 _{0.3}	–	0.4 _{0.1}	–	0.1 _{0.0}	–	331.8
RDV-Det*	1.8 _{0.5}	–	0.6 _{0.4}	–	0.3 _{0.1}	–	0.1 _{0.0}	–	146.6
Faster R-CNN	73.3 _{4.9}	63.2 _{5.4}	48.4 _{5.8}	42.1 _{5.3}	43.5 _{6.3}	37.6 _{5.4}	25.9 _{4.4}	22.6 _{3.9}	23.4
Cascade R-CNN	69.5 _{5.1}	59.6 _{5.6}	44.6 _{6.1}	38.5 _{5.6}	39.5 _{6.6}	33.6 _{5.6}	21.6 _{4.6}	18.7 _{4.1}	20.6
SSD	74.6 _{4.8}	64.5 _{5.3}	50.2 _{5.7}	43.7 _{5.2}	45.4 _{6.1}	39.4 _{5.2}	27.1 _{4.3}	23.8 _{3.9}	82.4
Vit-Det	86.5 _{2.8}	73.6 _{2.8}	52.2 _{4.6}	45.4 _{3.8}	48.1 _{4.8}	41.2 _{3.8}	30.2 _{3.9}	26.8 _{3.5}	16.8
Deformable-DETR	75.4 _{4.7}	65.0 _{5.2}	51.1 _{5.7}	44.5 _{5.1}	46.3 _{6.2}	40.1 _{5.2}	27.5 _{4.6}	24.0 _{4.1}	24.5
RT-DETR	91.6 _{0.9}	81.0 _{1.6}	58.9 _{4.7}	52.8 _{3.3}	56.8 _{2.4}	50.6 _{2.2}	33.0 _{3.8}	29.6 _{3.3}	66.3
YOLOv10	88.4 _{1.3}	80.7 _{2.2}	59.4 _{3.4}	<u>54.9</u> _{2.5}	54.6 _{3.2}	50.2 _{3.1}	34.3 _{4.1}	<u>31.8</u> _{3.5}	200.3
YOLOv11	88.2 _{1.4}	80.0 _{2.5}	59.1 _{3.2}	54.4 _{2.1}	55.6 _{3.6}	51.0 _{3.5}	34.1 _{3.7}	31.5 _{3.3}	185.2
YOLOv12	88.8 _{1.1}	80.4 _{1.9}	<u>59.9</u> _{3.1}	54.8 _{2.2}	54.5 _{2.1}	49.9 _{1.9}	<u>34.3</u> _{3.8}	31.5 _{3.2}	<u>204.1</u>
TAPIR	76.1 _{4.5}	65.8 _{4.8}	52.3 _{5.1}	45.6 _{4.6}	47.1 _{5.4}	40.5 _{4.7}	28.4 _{4.7}	24.6 _{4.2}	10.6
MCIT-IG	77.4 _{4.4}	67.2 _{4.6}	53.6 _{4.9}	46.9 _{4.3}	48.4 _{5.1}	41.8 _{4.5}	29.6 _{4.5}	26.0 _{4.0}	16.0
TDnet (Ours)	<u>89.9</u> _{1.3}	<u>81.0</u> _{2.0}	61.7 _{2.9}	56.3 _{2.1}	<u>55.7</u> _{2.4}	<u>50.8</u> _{2.7}	36.1 _{3.4}	33.1 _{3.1}	126.6

* Weakly-supervised methods.

Thanks for your attention

:)

Yiliang Chen