



# Unsupervised Surgical Instrument Segmentation via Anchor Generation and Semantic Diffusion

Daochang Liu\*, Yuhui Wei\*, Tingting Jiang, Yizhou Wang, Rulin Miao, Fei Shan, and Ziyu Li



北京大学  
PEKING UNIVERSITY



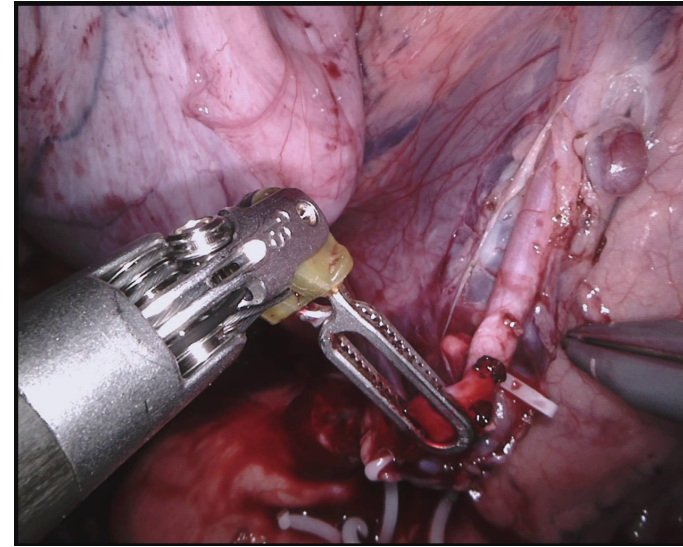
北京大学 肿瘤医院  
BEIJING CANCER HOSPITAL



# Introduction

---

## Instrument Segmentation In Surgical Video



U-Net



# Introduction

---

## Existing Fully Supervised Methods

Very Expensive



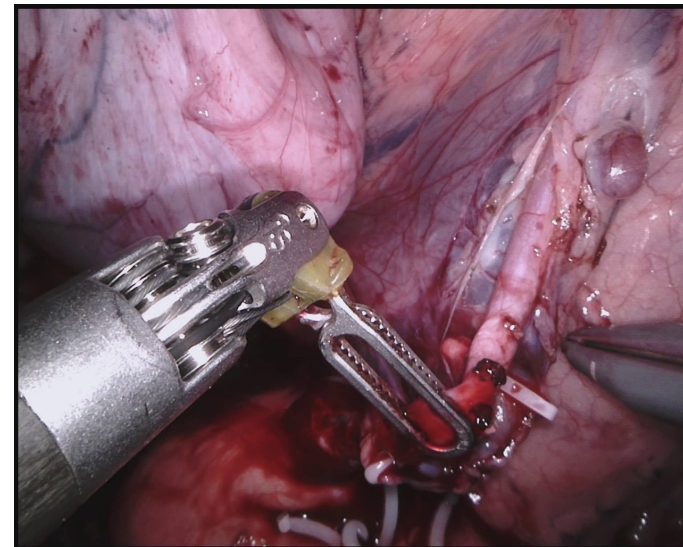
...



...

Manual  
Annotation

## Instrument Segmentation In Surgical Video



# Introduction

## Existing Fully Supervised Methods

Very Expensive



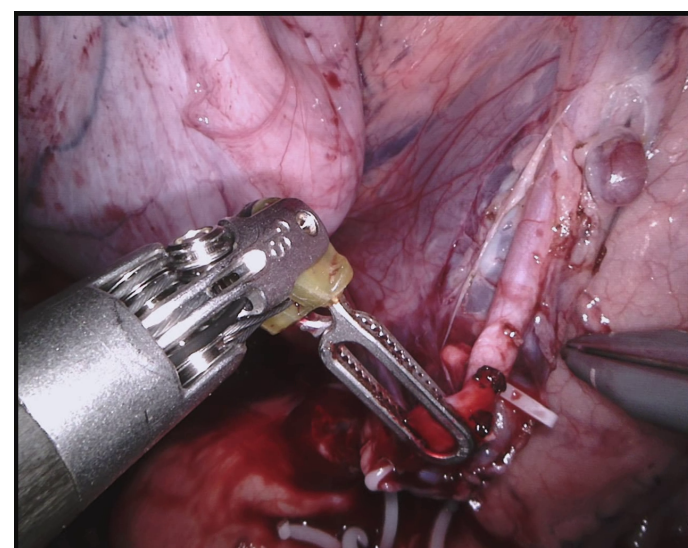
...



...

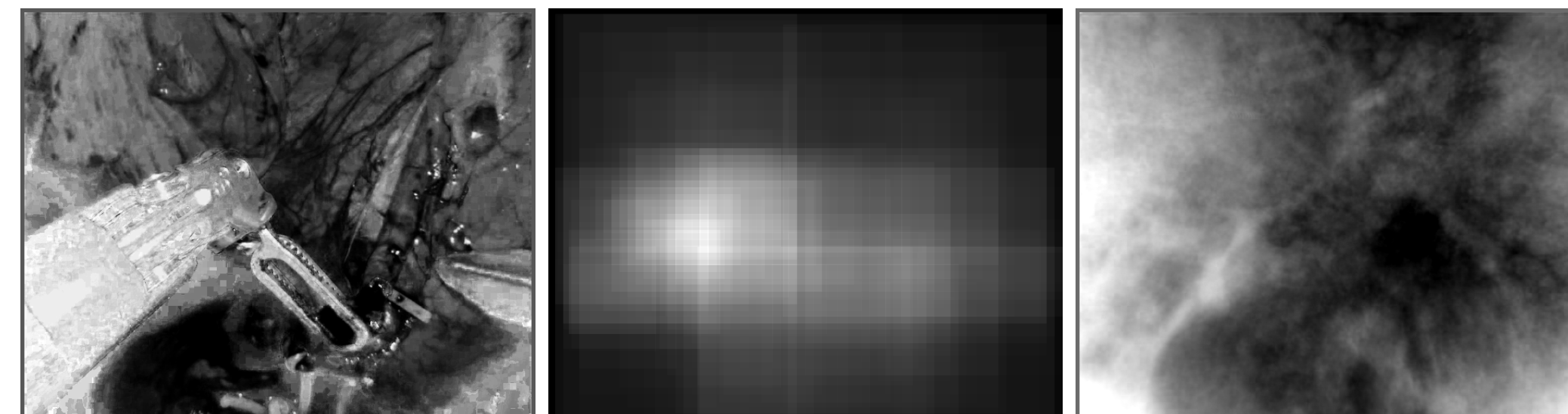
Manual Annotation

## Instrument Segmentation In Surgical Video



## Our Proposed Unsupervised Method

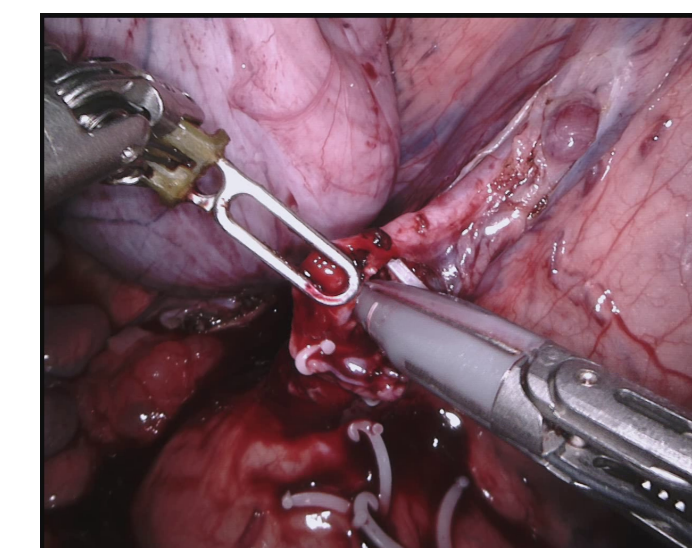
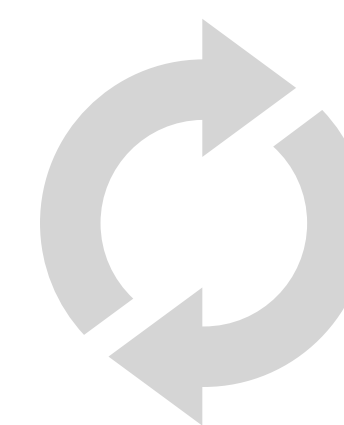
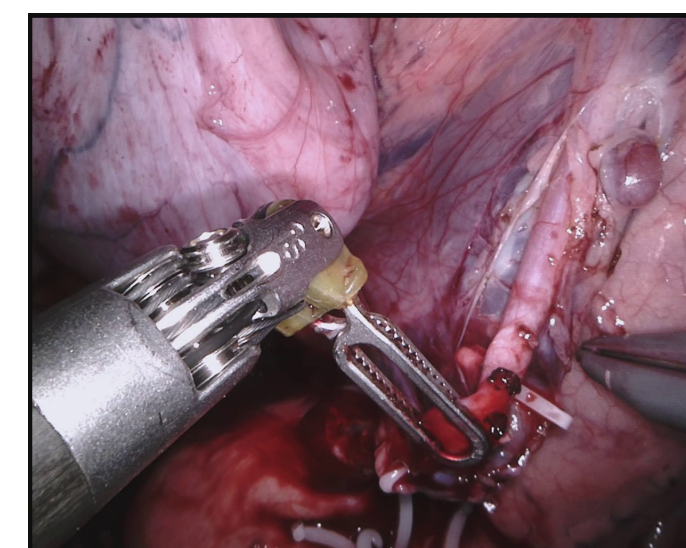
More Affordable 



Anchor Generation

Coarse Cues (Prior Knowledge)

Semantic Diffusion

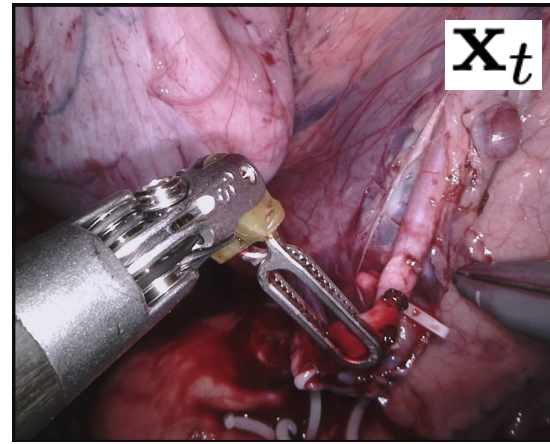


Inter-image Correlation

# Method

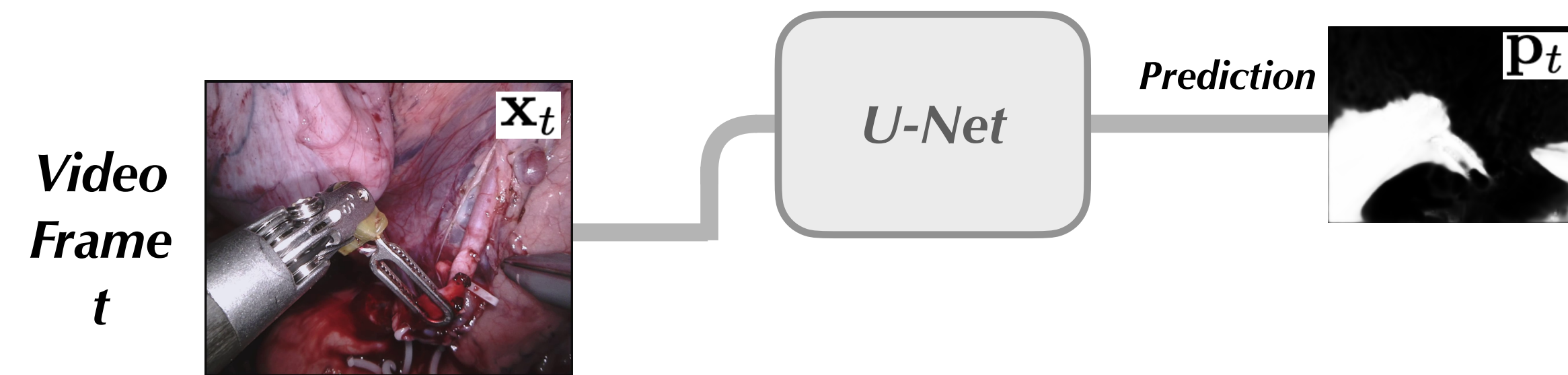
---

*Video*  
*Frame*  
*t*



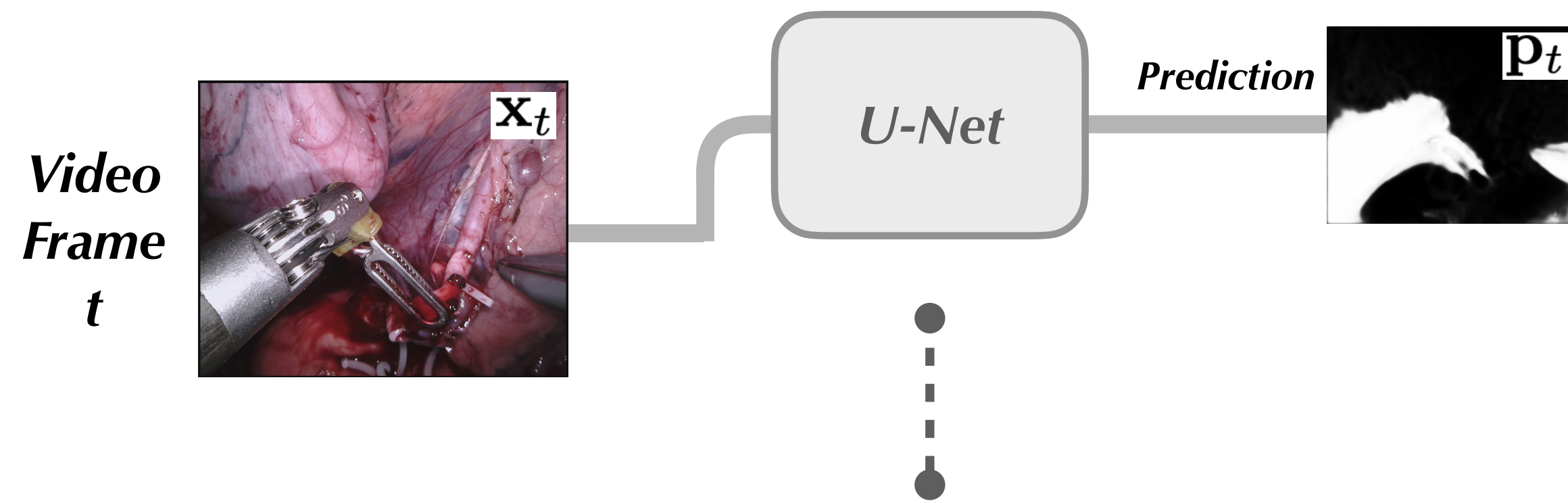
# Method

---



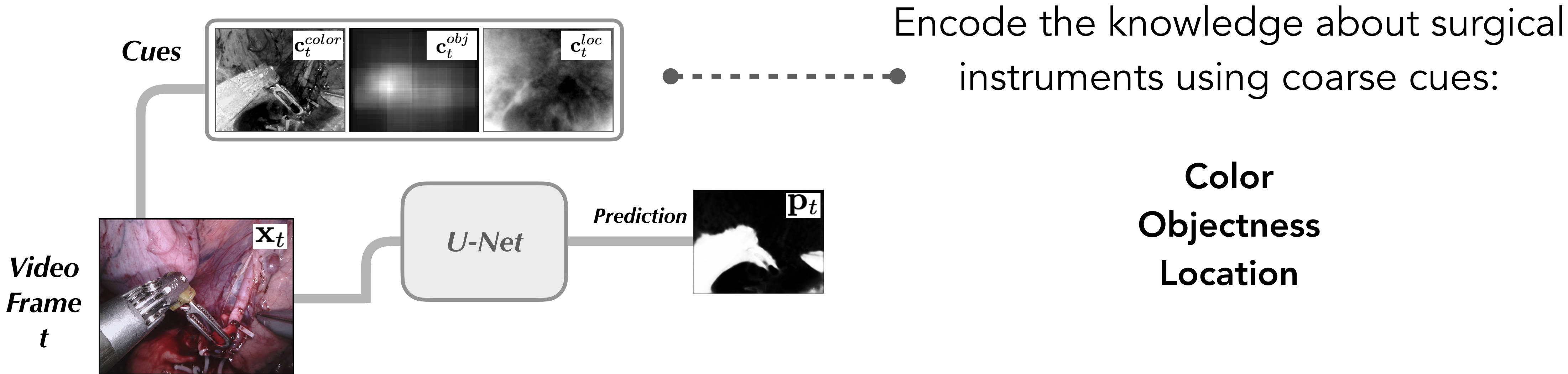
# Method

---



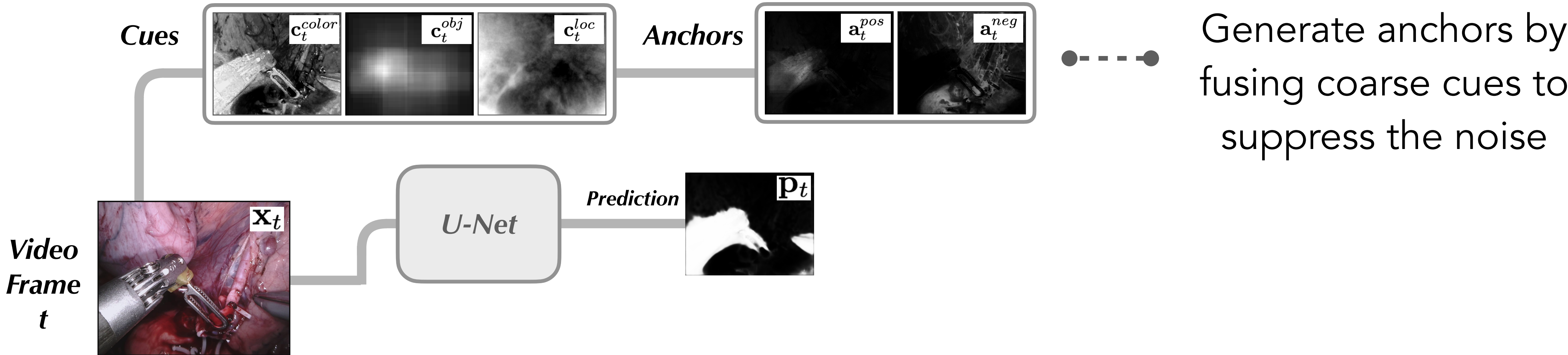
How to provide training supervision without human annotation?

# Method

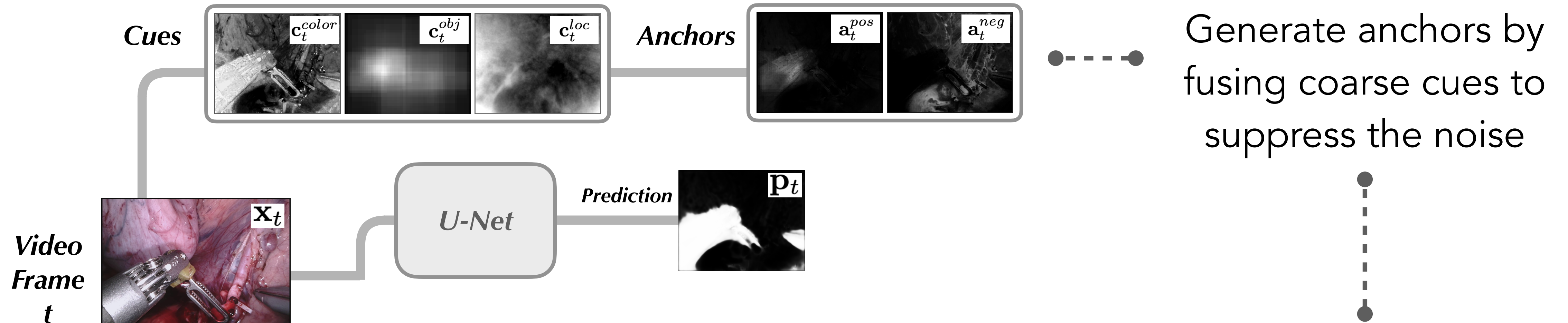




# Method



# Method



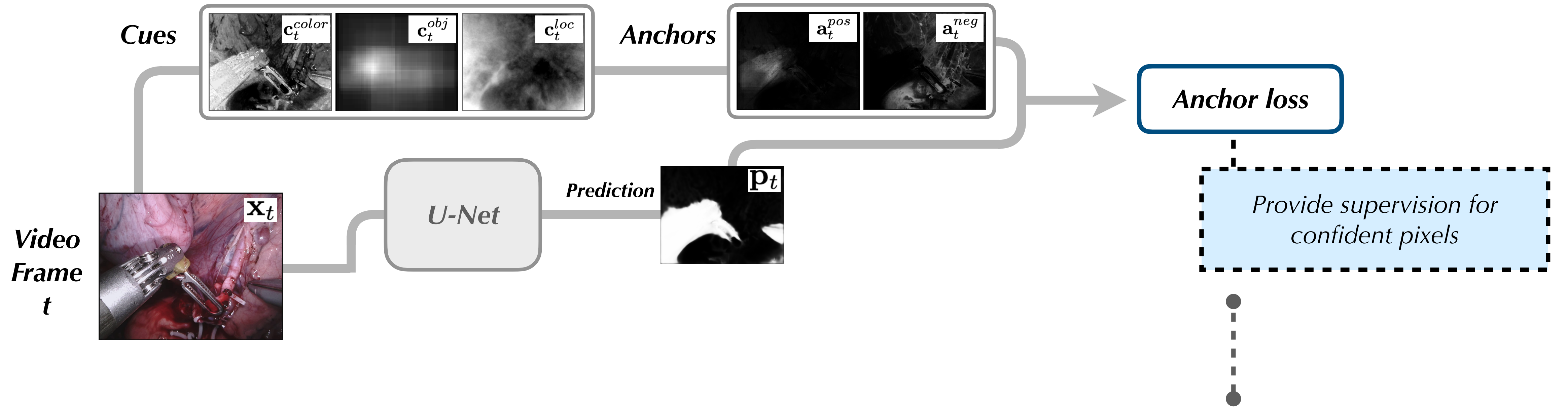
**Positive anchor** captures the confident instrument regions that satisfy all the cues

$$\mathbf{a}_t^{pos} = \mathbf{c}_t^{color} \mathbf{c}_t^{obj} \mathbf{c}_t^{loc}$$

**Negative anchor** captures the confident background regions that satisfy none of the cues

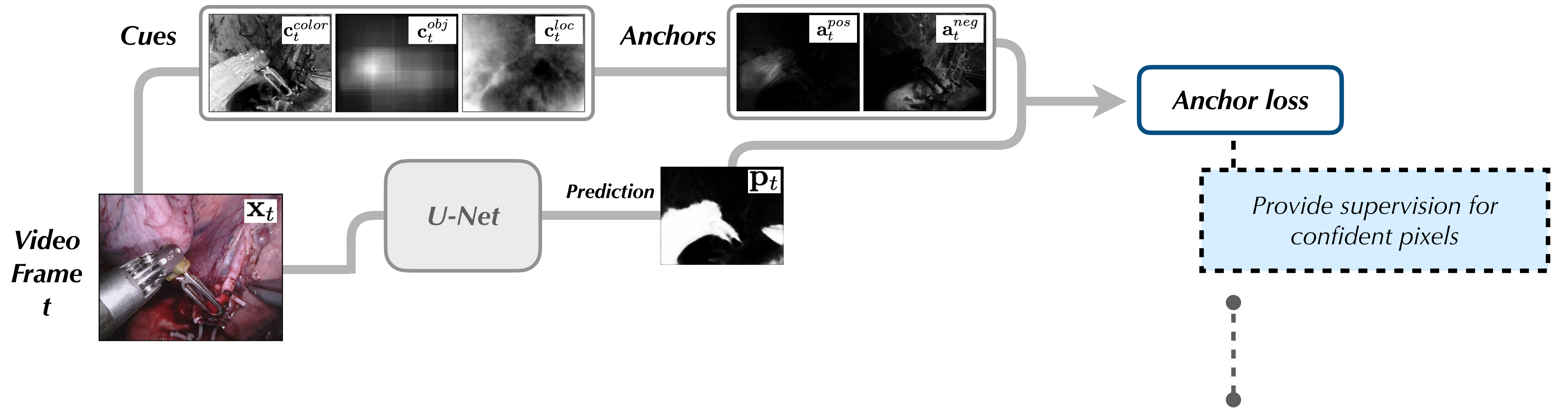
$$\mathbf{a}_t^{neg} = (1 - \mathbf{c}_t^{color})(1 - \mathbf{c}_t^{obj})(1 - \mathbf{c}_t^{loc})$$

# Method



Anchors are used as pseudo labels to supervise confident pixels

# Method

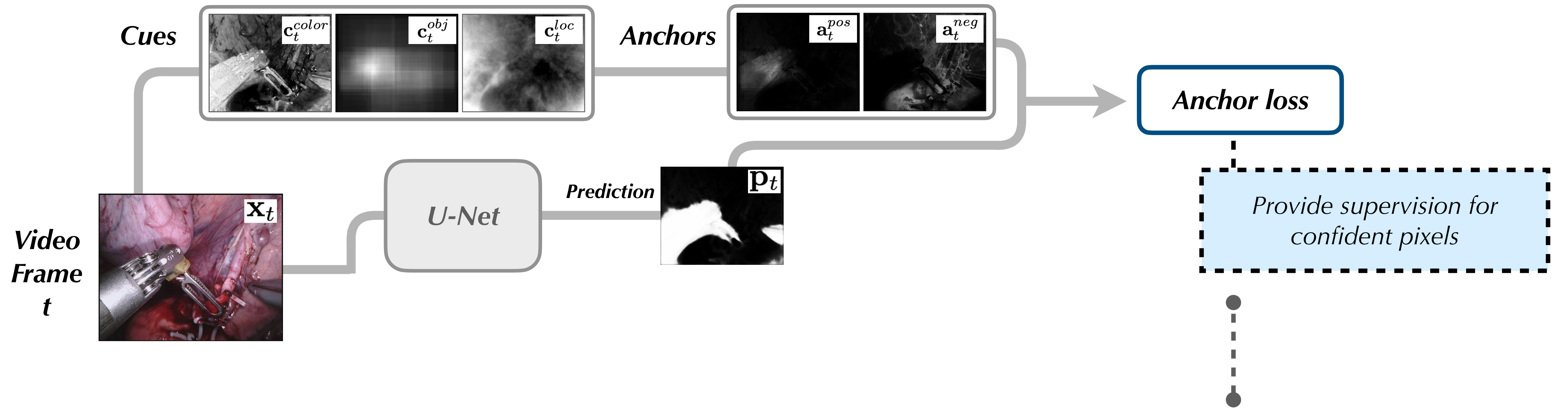


## Anchor loss

$$\mathcal{L}_{anc}(\mathbf{x}_t) = \frac{1}{HW} \sum_i -\mathbf{a}_{t,i}^{pos} \mathbf{p}_{t,i} - \mathbf{a}_{t,i}^{neg} (1 - \mathbf{p}_{t,i})$$

●-----● Anchors are used as pseudo labels to supervise confident pixels

# Method



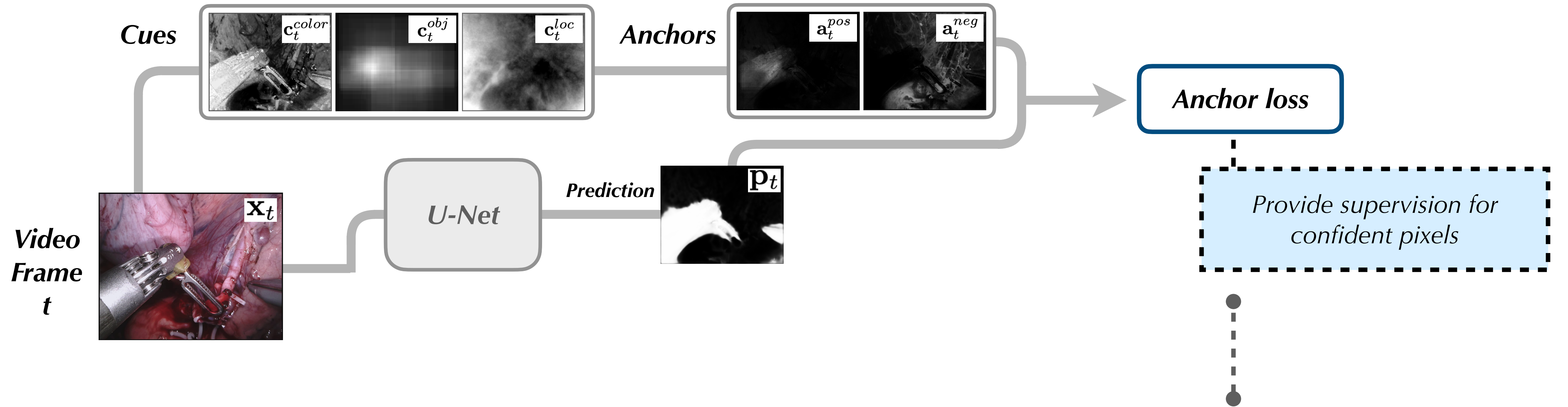
## Anchor loss

$$\mathcal{L}_{anc}(\mathbf{x}_t) = \frac{1}{HW} \sum_i -\mathbf{a}_{t,i}^{pos} \mathbf{p}_{t,i} - \mathbf{a}_{t,i}^{neg} (1 - \mathbf{p}_{t,i})$$

Anchors are used as pseudo labels to supervise confident pixels

Encourage activation on the positive anchor  
Inhibit activation on the negative anchor

# Method



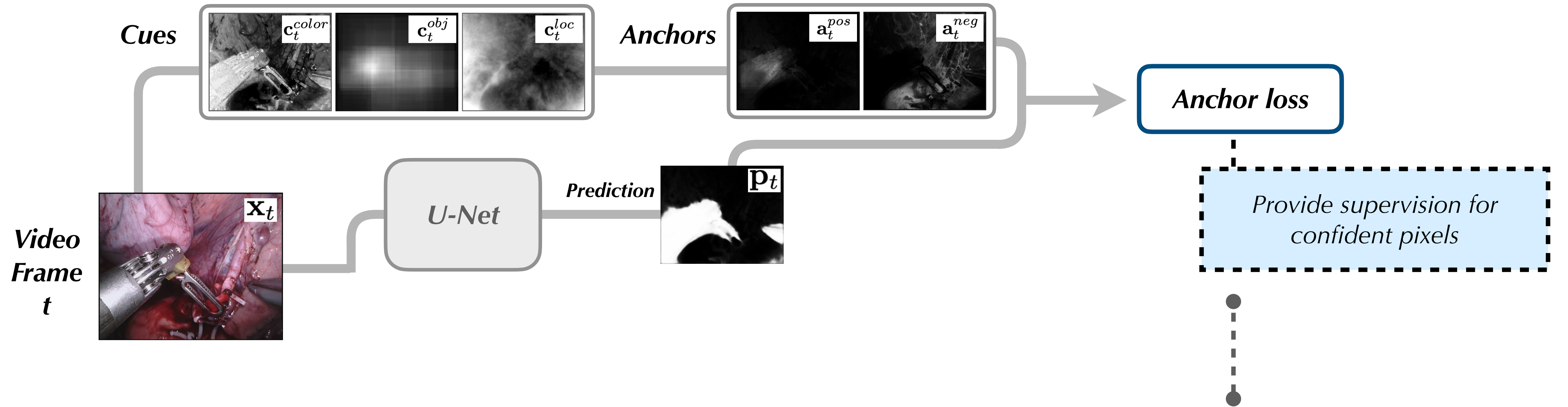
## Anchor loss

$$\mathcal{L}_{anc}(\mathbf{x}_t) = \frac{1}{HW} \sum_i -\mathbf{a}_{t,i}^{pos} \mathbf{p}_{t,i} - \mathbf{a}_{t,i}^{neg} (1 - \mathbf{p}_{t,i})$$

Anchors are used as pseudo labels to supervise confident pixels

How to supervise the remaining ambiguous pixels outside the anchors?

# Method



## Anchor loss

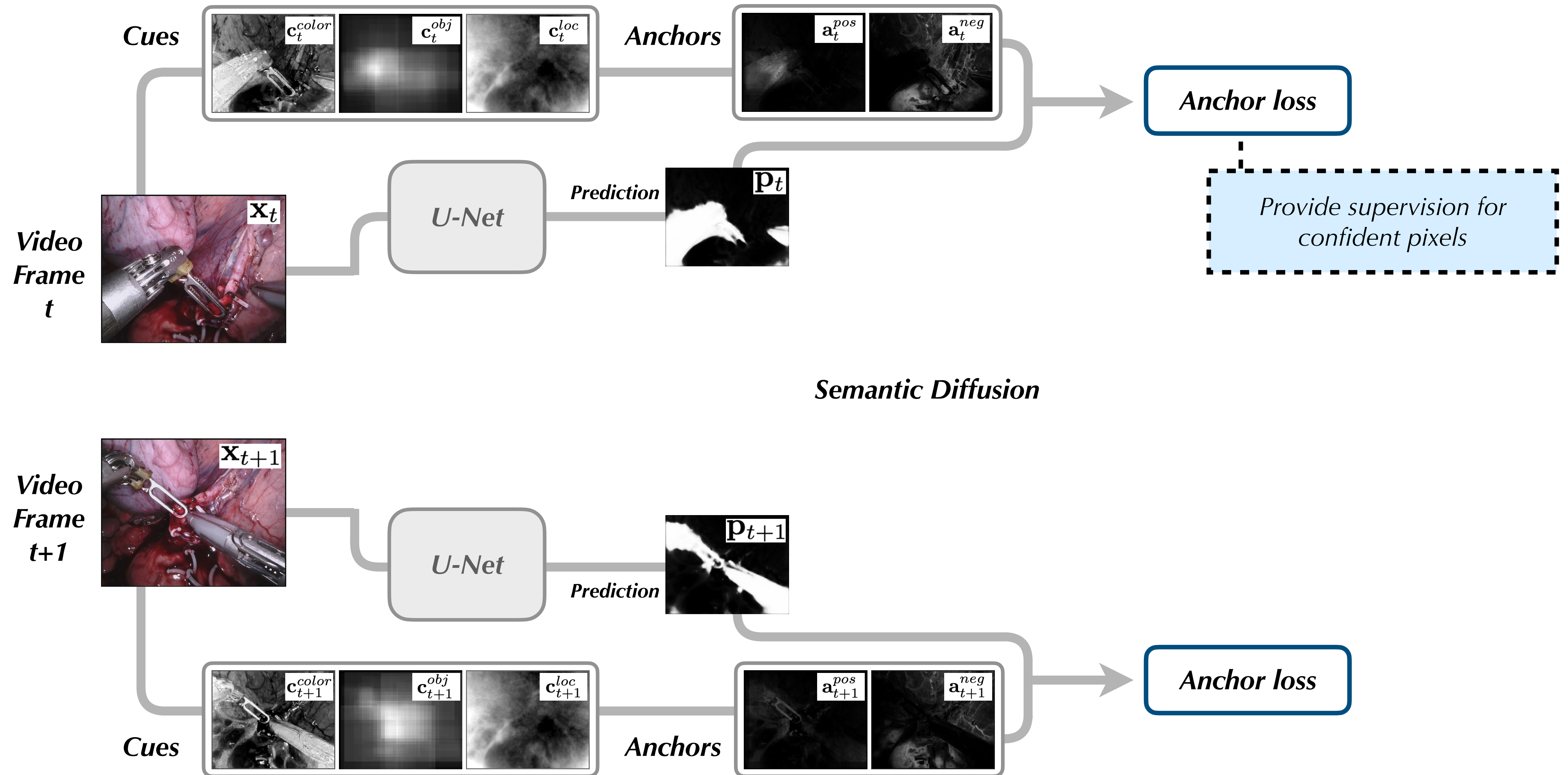
$$\mathcal{L}_{anc}(\mathbf{x}_t) = \frac{1}{HW} \sum_i -\mathbf{a}_{t,i}^{pos} \mathbf{p}_{t,i} - \mathbf{a}_{t,i}^{neg} (1 - \mathbf{p}_{t,i})$$

Anchors are used as pseudo labels to supervise confident pixels

How to supervise the remaining ambiguous pixels outside the anchors?

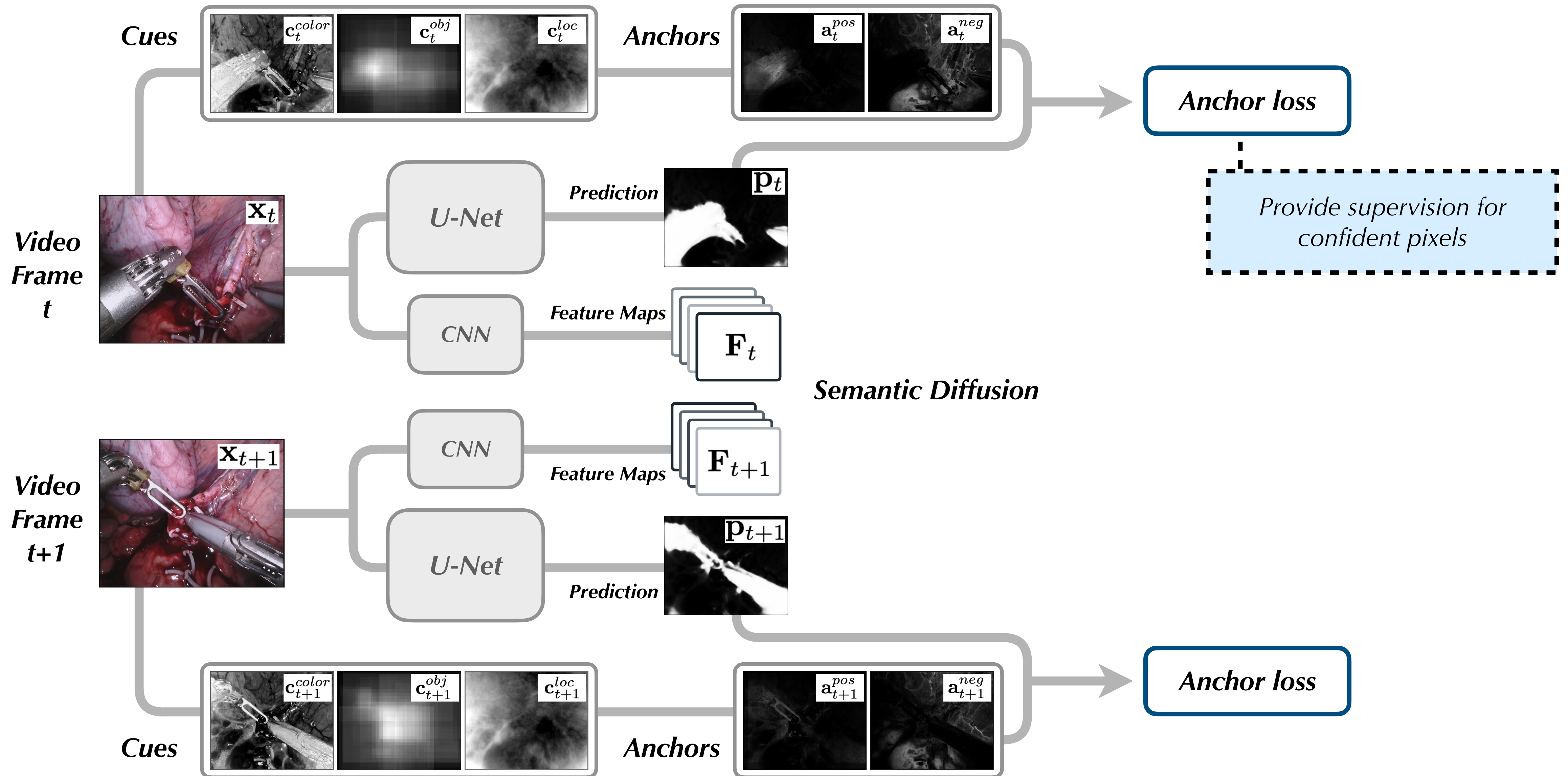
## Semantic Diffusion

# Method

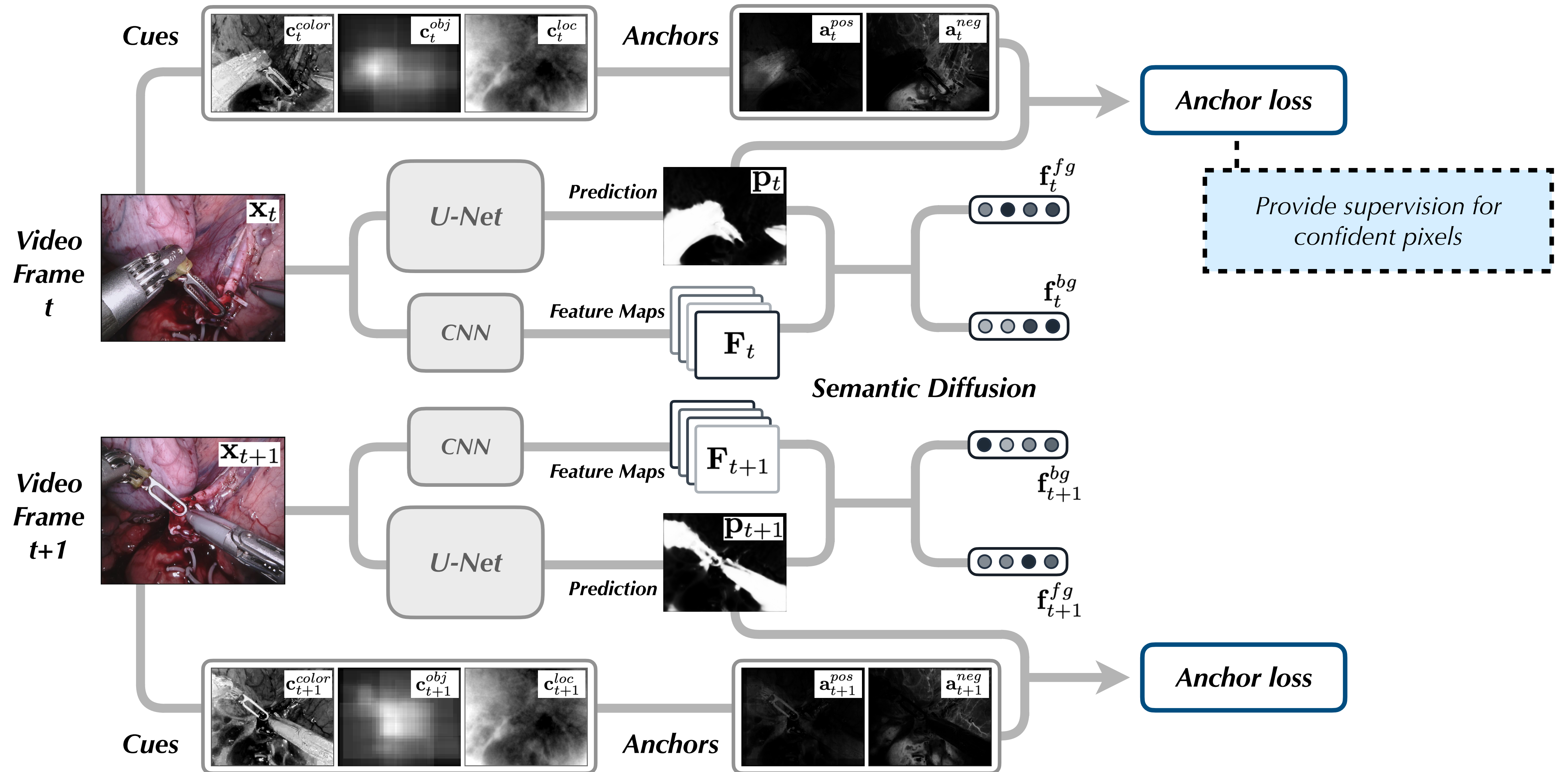




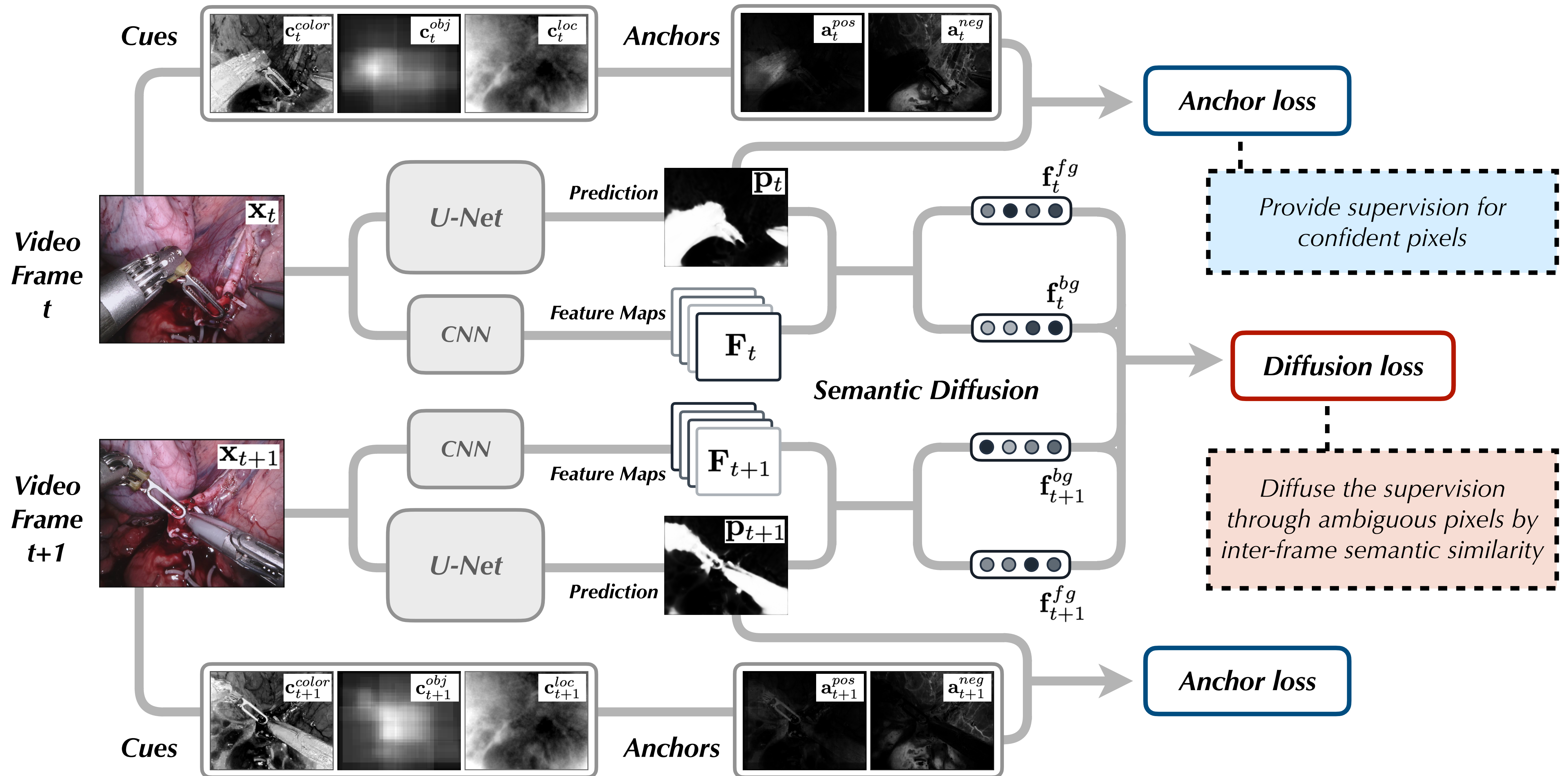
# Method



# Method



# Method



# Method

---

## Diffusion loss

$$\mathcal{L}_{dif}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg}) - 2\phi(\mathbf{f}_t^{fg}, \mathbf{f}_{t+1}^{fg}) + m^{fg}, 0)$$

# Method

---

## Diffusion loss

$$\mathcal{L}_{diff}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\underbrace{\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg})}_{\text{Semantic features of the foreground and the background regions in the two frames}} + \underbrace{\phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg})}_{\text{Semantic features of the foreground and the background regions in the two frames}} - 2\underbrace{\phi(\mathbf{f}_t^{fg}, \mathbf{f}_{t+1}^{fg})}_{\text{Adjacent video frames}} + \underbrace{m^{fg}}_{\text{Margin hyperparameter}}, 0)$$

Adjacent video frames

Semantic features of the foreground and the background regions in the two frames

Margin hyperparameter

# Method

---

## Diffusion loss

$$\mathcal{L}_{dif}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\underbrace{\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg})}_{\text{Minimize intra-frame instrument-background similarities}} - \underbrace{2\phi(\mathbf{f}_t^{fg}, \mathbf{f}_{t+1}^{fg})}_{\text{Maximize inter-frame instrument-instrument similarity}} + m^{fg}, 0)$$

Minimize intra-frame instrument-background similarities

Maximize inter-frame instrument-instrument similarity

# Method

---

## Diffusion loss

$$\mathcal{L}_{dif}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg}) - 2\phi(\mathbf{f}_t^{fg}, \mathbf{f}_{t+1}^{fg}) + m^{fg}, 0)$$

$$\mathcal{L}_{dif}^{bg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg}) - \underline{2\phi(\mathbf{f}_t^{bg}, \mathbf{f}_{t+1}^{bg})} + m^{bg}, 0)$$



Similarly, maximize the inter-frame background-background similarity

# Method

---

## Diffusion loss

$$\mathcal{L}_{dif}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg}) - 2\phi(\mathbf{f}_t^{fg}, \mathbf{f}_{t+1}^{fg}) + m^{fg}, 0)$$

$$\mathcal{L}_{dif}^{bg}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \max(\phi(\mathbf{f}_t^{fg}, \mathbf{f}_t^{bg}) + \phi(\mathbf{f}_{t+1}^{fg}, \mathbf{f}_{t+1}^{bg}) - 2\phi(\mathbf{f}_t^{bg}, \mathbf{f}_{t+1}^{bg}) + m^{bg}, 0)$$

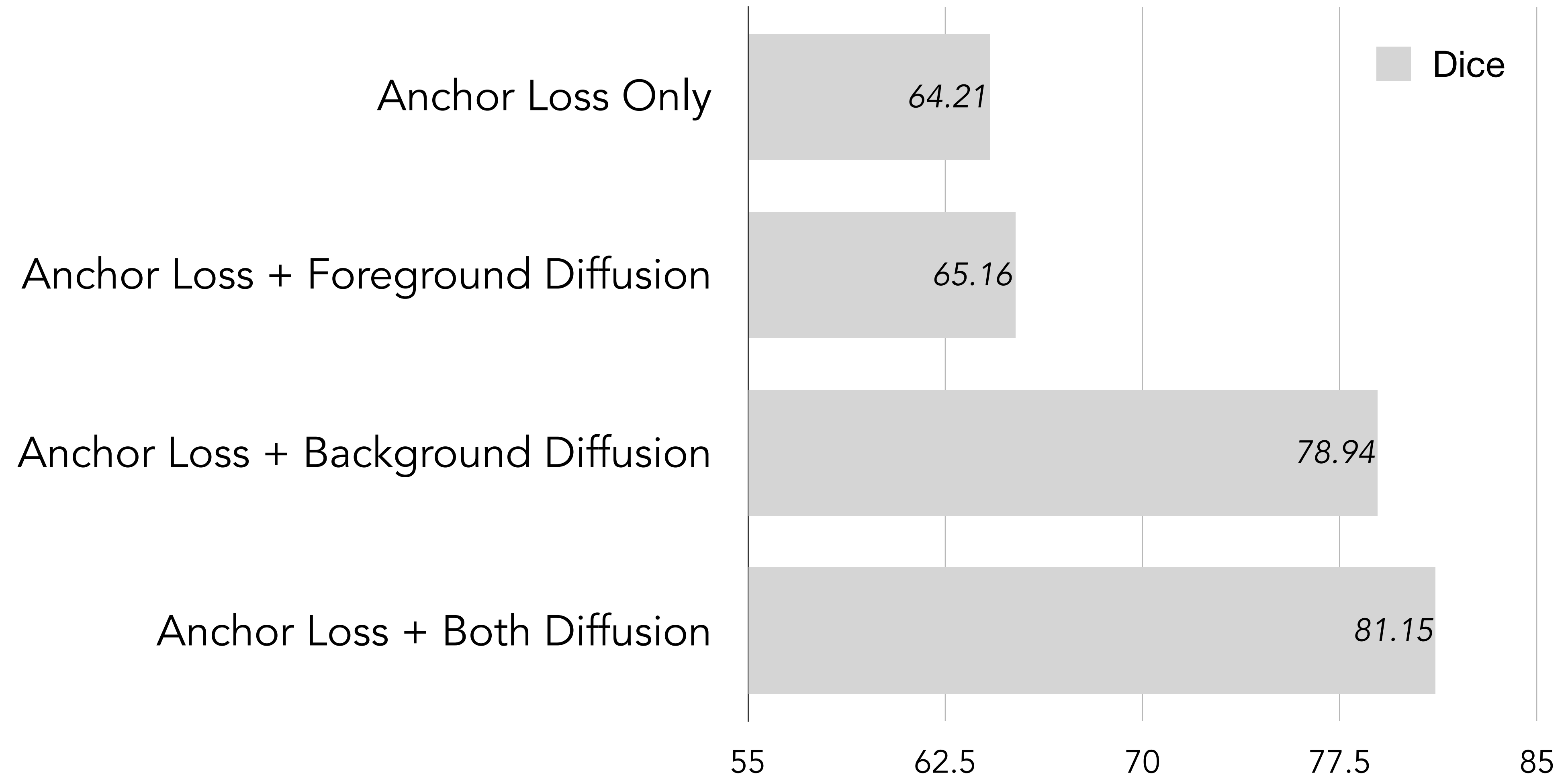
**Full loss = Anchor loss + Diffusion loss**

$$\mathcal{L}_{full}(\mathbf{x}_t, \mathbf{x}_{t+1}) = \mathcal{L}_{anc}(\mathbf{x}_t) + \mathcal{L}_{anc}(\mathbf{x}_{t+1}) + \mathcal{L}_{dif}^{fg}(\mathbf{x}_t, \mathbf{x}_{t+1}) + \mathcal{L}_{dif}^{bg}(\mathbf{x}_t, \mathbf{x}_{t+1})$$



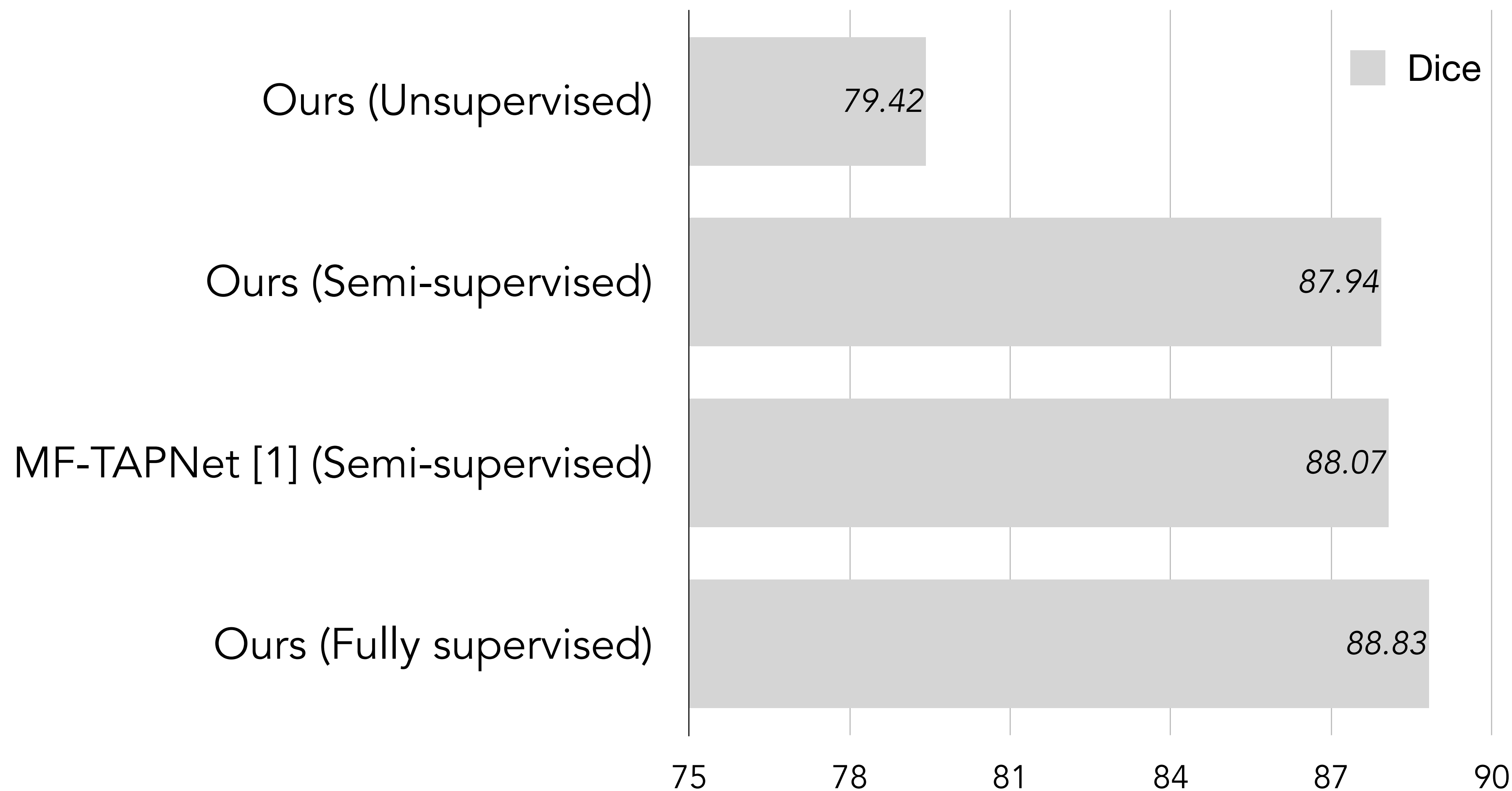
# Results on EndoVis 2017

---



# Extend to Semi-Supervised and Fully Supervised Settings

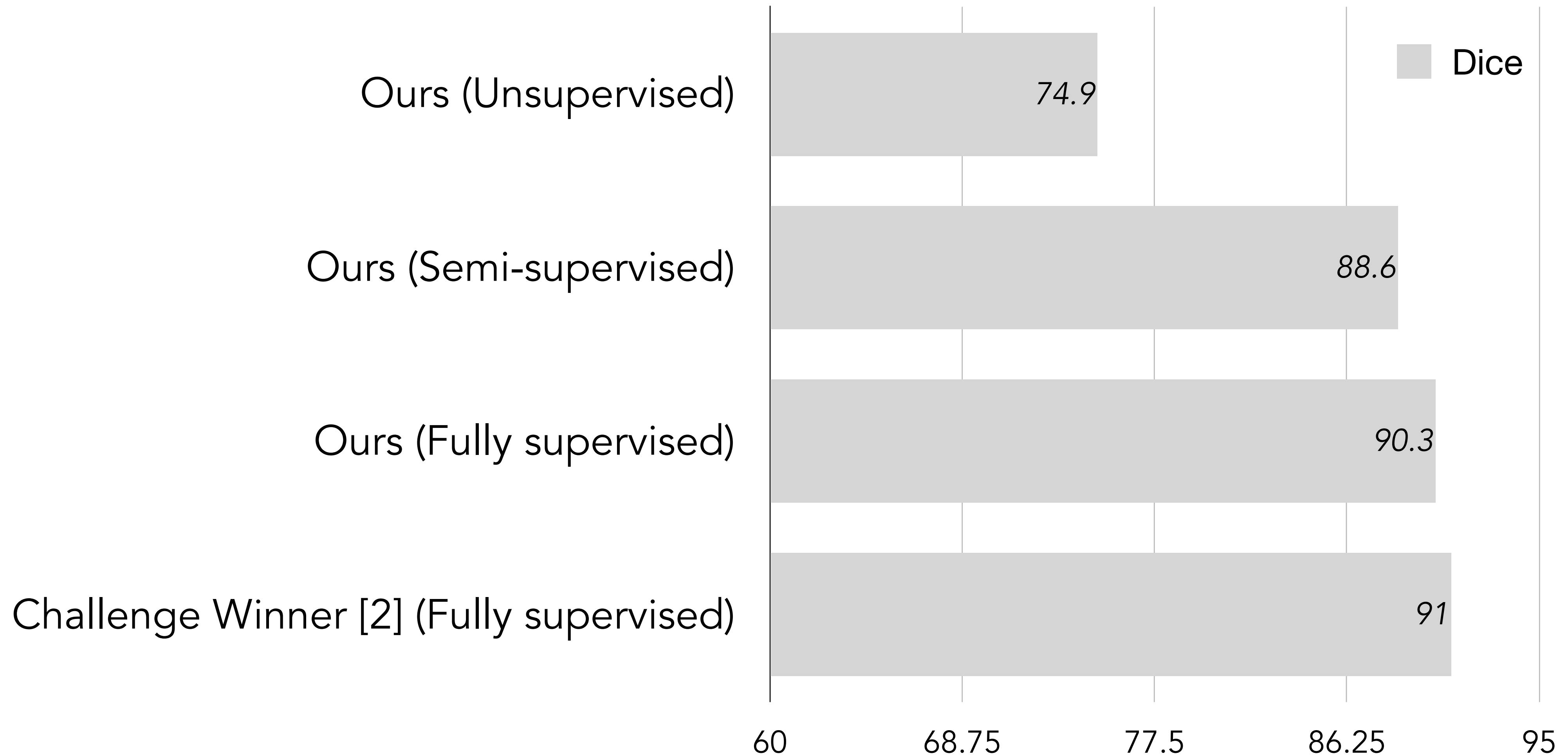
---



**We only use vanilla U-Net**

# Extend to Skin Lesion Segmentation on ISIC 2016

---

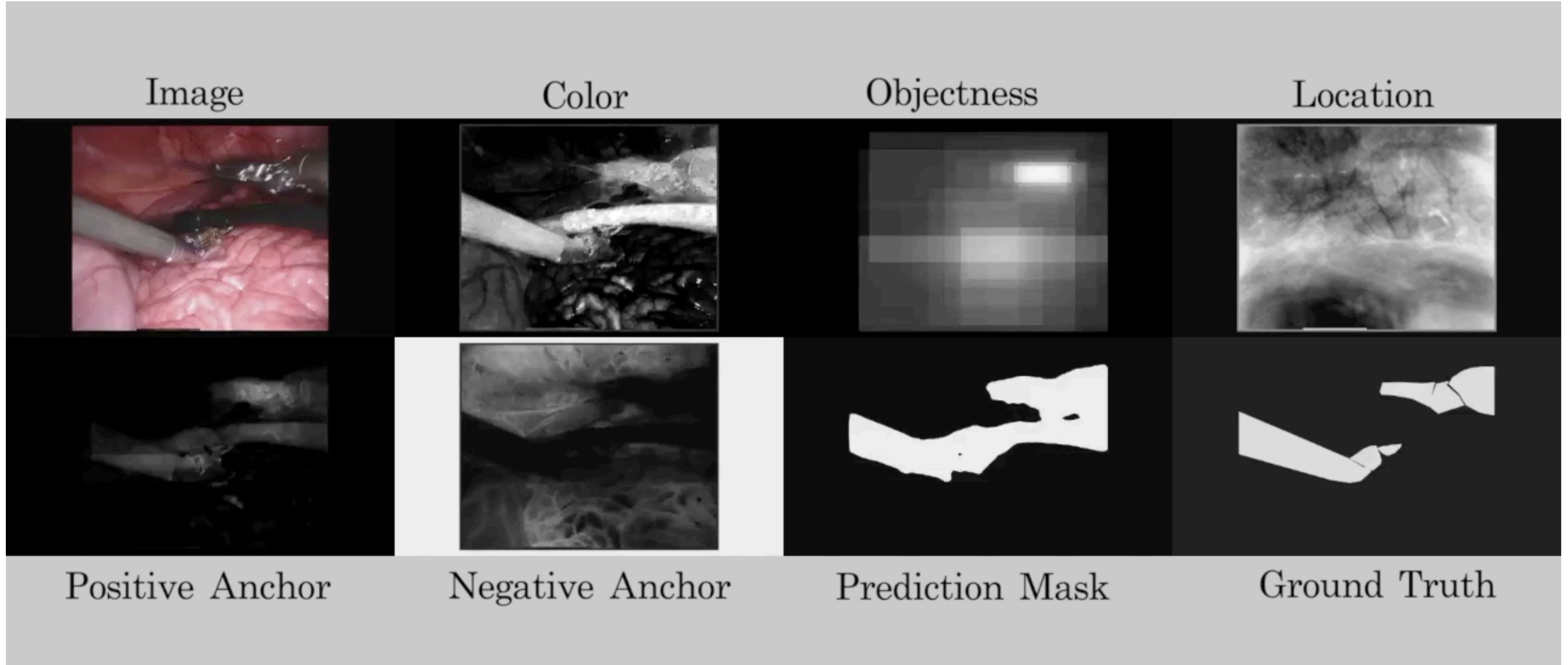


**We only use vanilla U-Net**

[2] Gutman, D., Codella, N.C., Celebi, E., Helba, B., Marchetti, M., Mishra, N., Halpern, A.: Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (ISBI) 2016, hosted by the international skin imaging collaboration (ISIC). arXiv:1605.01397 (2016)

# Video Demo

---



# Conclusion

---

- **Unsupervised** Surgical Instrument Segmentation
- No human annotation needed
- **Anchor generation** to supervise confident pixels
- **Semantic diffusion** to supervise ambiguous pixels
- Promising results on **EndoVis 2017** and **ISIC 2016**
- **Codes** at <https://github.com/Finspire13/AGSD-Surgical-Instrument-Segmentation>

***Thank you!***

***Welcome to Discuss Online.***

ORAL SESSION WEDS-AM-4

6:30 PM - 7:00 PM CST on Wednesday, 7 October

POSTER SESSION 5

5:00 PM - 6:30 PM CST on Wednesday, 7 October

