

YouTube-VOS: Sequence-to-Sequence Video Object Segmentation  
ECCV 2018  
[arXiv: Sep 3, 2018]

YouTube-VOS: A Large-Scale Video Object Segmentation Benchmark  
[arXiv: Sep 6, 2018]

---

Presented By:  
Jyoti Kini

# Video Object Segmentation

---

- ❑ Segmenting a particular object instance throughout the entire video sequence given only the object mask on the first frame



# YouTube-VOS Dataset

Scale	YouTube-VOS (Ours)	YouTube-VOS (Ours)
Videos	<b>3,252</b>	<b>4,453</b>
Categories	<b>78</b>	<b>94</b>
Objects	<b>6,048</b>	<b>7,755</b>
Annotations	<b>133,886</b>	<b>197,272</b>
Duration	<b>217.21</b>	<b>334.81</b>

[Partial dataset used for model implementation] [Final dataset]

- ❑ Human annotators select up to five objects per video clip and annotate every five frames in a 30fps frame rate, resulting in 6fps sampling rate

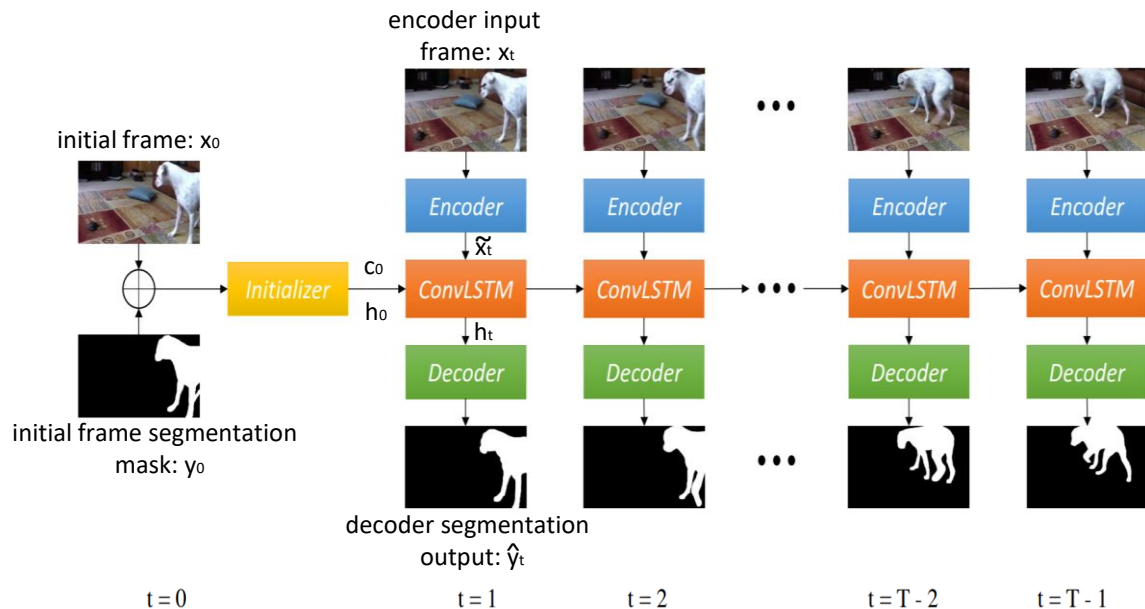
Videos - 4,453

Training - 3,471

Validate – 474  
(26 unseen)

Testing – 508  
(26 unseen)

# Algorithm



$$c_0, h_0 = \text{Initializer}(x_0, y_0)$$

$$\tilde{x}_t = \text{Encoder}(x_t)$$

$$c_t, h_t = \text{ConvLSTM}(\tilde{x}_t, c_{t-1}, h_{t-1})$$

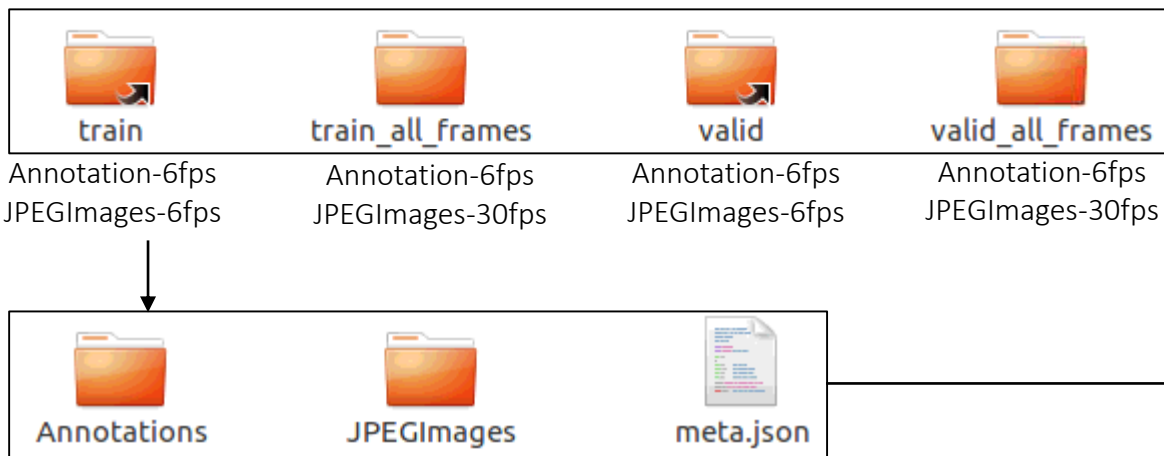
$$\hat{y}_t = \text{Decoder}(h_t)$$

$$\mathcal{L} = -(y_t \log(\hat{y}_t)) + ((1 - y_t) \log(1 - \hat{y}_t))$$

# Implementation Details

## □ Data Loader:

- Dataset directories:



```
"0043f083b5": {  
  "objects": {  
    "1": {  
      "category": "bus",  
      "frames": [  
        "0000",  
        "0005",  
        "0010",  
        "0015",  
        "0020",  
        "0025",  
        "0030",  
        "0035",  
        "0040",  
        "0045",  
        "0050",  
        "0055",  
        "0060",  
        "0065",  
        "0070",  
        "0075",  
        "0080",  
        "0085",  
        "0090",  
        "0095"  
      ]  
    },  
    "2": {  
      "category": "sedan",  
      "frames": [  
        "0010",  
        "0015",  
        "0020",  
        ]  
    }  
  }  
}
```

# Implementation Details

---

## □ Data Loader:

- JSON:
  - Not all objects appear from the 1<sup>st</sup> frame onwards
  - Colors to binary masks (PIL/OpenCV)
  - Handle occlusion

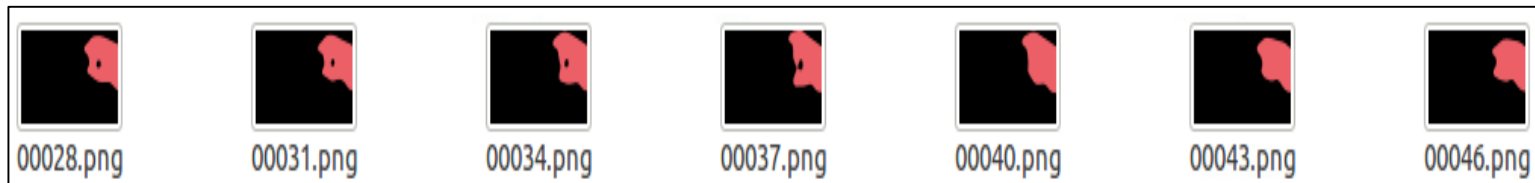
```
"0043f083b5": {  
  "objects": {  
    "1": {  
      "category": "bus",  
      "frames": [  
        "0000",  
        "0005",  
        "0010",  
        "0015",  
        "0020",  
        "0025",  
        "0030",  
        "0035",  
        "0040",  
        "0045",  
        "0050",  
        "0055",  
        "0060",  
        "0065",  
        "0070",  
        "0075",  
        "0080",  
        "0085",  
        "0090",  
        "0095"  
      ]  
    },  
    "2": {  
      "category": "sedan",  
      "frames": [  
        "0010",  
        "0015",  
        "0020",  
      ]  
    }  
  }  
}
```

# Implementation Details

---

## □ Data Loader:

- CPU threads
- Random selection of single object instance (PIL/OpenCV)

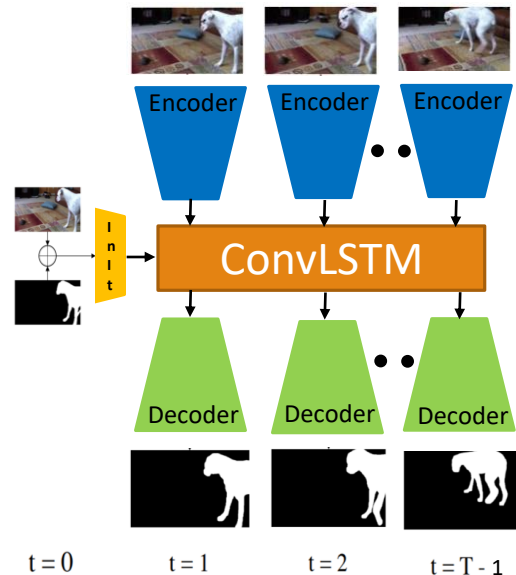


- Select time frames T (5 ~ 11) from a random training video sequence
- Resize original RGB frames and annotations to 256×448

# Implementation Details

## □ Model structure:

Component	Layer	Filters	Filter size	Initialization	Activation
Initializer	VGG-16 (all convolution)				
	VGG-16 (first fully-connected)		1x1		
Encoder	↳ Convolution	512	1x1	Xavier	Relu
	↳ Convolution	512	1x1	Xavier	Relu
Encoder	VGG-16 (all convolution)				
	VGG-16 (first fully-connected)		1x1		
ConvLSTM	Convolution	512	1x1	Xavier	Relu
	ConvLSTM			Xavier	Sigmoid - gate outputs Relu - state outputs
Decoder	ConvLSTM	512	3x3	Forget bias=1	
	Deconvolution	512	5x5	Xavier	
	Deconvolution	256	5x5	Xavier	
	Deconvolution	128	5x5	Xavier	
	Deconvolution	64	5x5	Xavier	
	Deconvolution	64	5x5	Xavier	
	Convolution	1	5x5	Xavier	Sigmoid





# Implementation Details

---

## □ Training:

- Primarily, trained on 6fps, followed by 30fps as the loss stabilizes
- Loss for frames without ground-truth is set to 0
- Model converges in 80 epochs
- Initial learning rate is set to  $10^{-5}$

# Implementation Details

## □ Evaluation:

- Submission access URL

(<https://competitions.codalab.org/competitions/19544>)

- Submission format



## Input:

```
<submission>.zip
  |- Annotations
    |- <video_id>
      |- <frame_id>.png
      |- <frame_id>.png
    |- <video_id>
      |- <frame_id>.png
      |- <frame_id>.png
```

## Output:

```
Overall: 0.474763296366
J_seen: 0.548071686125
J_unseen: 0.419488520321
F_seen: 0.495602588958
F_unseen: 0.315890390059
```

# Segmentation Metrics

---

- ❑ **Region similarity:** Defined as the intersection-over-union between the estimated segmentation and the ground truth mask

$$\mathcal{J} = \frac{|M \cap G|}{|M \cup G|}$$

Region Similarity is the intersection-over-union between mask M and ground truth G

- ❑ **Contour accuracy:** Interprets the masks as a set of closed contours and computes the contour-based F-measure which is a function of precision and recall

$$\mathcal{F} = \frac{2P_c R_c}{P_c + R_c}$$

Contour Accuracy is the F-measure for the contour based precision and recall

# Questions

---