# STESA-GRU: Spatio Temporal Self Attended Gated Recurrent Units for Video Object Segmentation

Diego Valderrama & Mateo Rueda

# INTRODUCTION

# INTRODUCTION

- Separate objects from a video sequence background

- Beneficial applications
  - Video editing
  - Scene understanding
  - Autonomous vehicles

# INTRODUCTION

- Semi-supervised task

- Predict masks from the first frame ground-truth
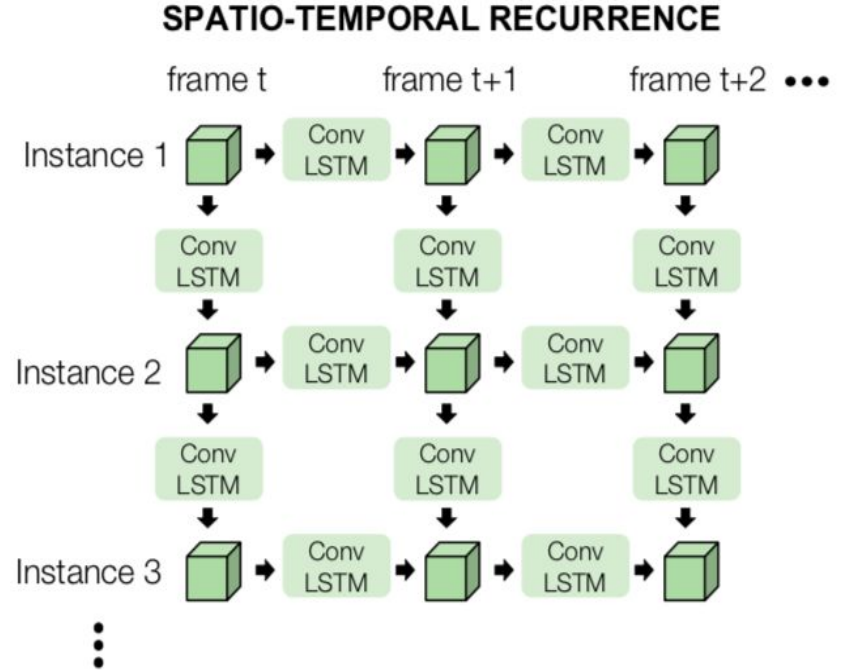
- Motion based or detection based approaches

# BASELINE

# RVOS

- Encoder - Decoder architecture

- Resnet 101 as backbone

- LSTM spatio-temporal

- Segmentation per instance

- Off-line approach

**SPATIO-TEMPORAL RECURRENCE**

# APPROACH

# APPROACH

- Conv GRU spatio-temporal

  - Two hidden states
  - Conv reset gate
  - Conv update gate
  - Conv "current"
  - Reset hidden

$$r_g = \sigma\left(W_{xr}x_t + W_{hr}h_{t-1} + b_r\right)$$

$$u_g = \sigma\left(W_{xu}x_t + W_{hu}h_{t-1} + b_u\right)$$

$$c = \sigma\left(W_{xc}x_t + b_c\right)$$

$$r_h = r_g * \left(hS_{t-1} + hT_{t-1}\right)$$

$$\tilde{h}_t = tanh\left(c + r_h\right)$$

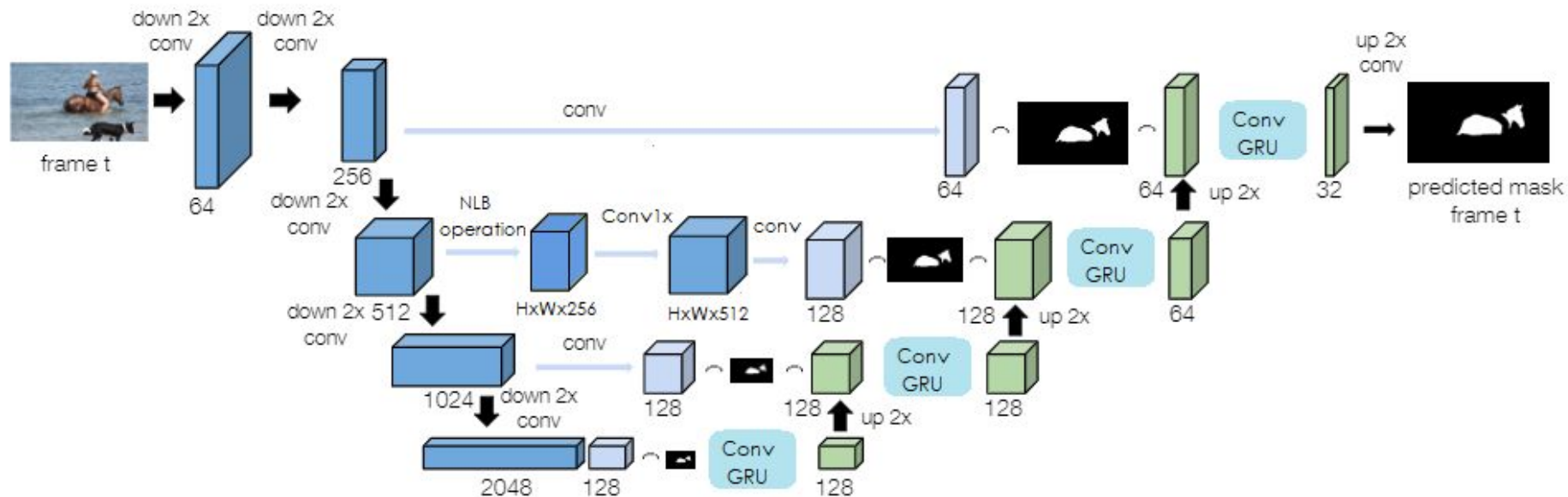$$h = (u_g * hS_{t-1}) + (1 - u_g) * \tilde{h}_t$$

# APPROACH

- Self- Attention

  - 2D Non Local Block Module after *res3*

  - Auxiliary Loss

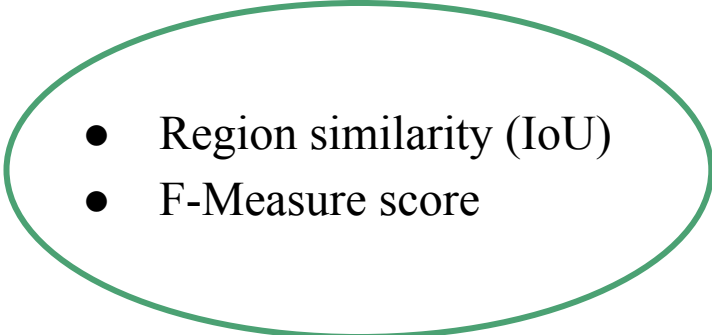$$Loss = mIoU_{NLB} + mIoU_{P}$$

# STESA-GRU

# EXPERIMENTS

# DAVIS 2017

# YOUTUBE VOS V1

- 60 videos - Train split
- 30 videos - Val split
- 30 videos - Test-dev

- 3471 videos - Train split

- Region similarity (IoU)
- F-Measure score

# **TRAINING DETAILS**

- Frames and annotations have been resized to 256x448

- Each mini-batch is composed of 4/3 videos and 5 consecutive frames

- 5/10 epochs in YouTube-VOS and 40/20 in DAVIS 2017

- Learning rate of $10^{-6}$ with Adam optimizer
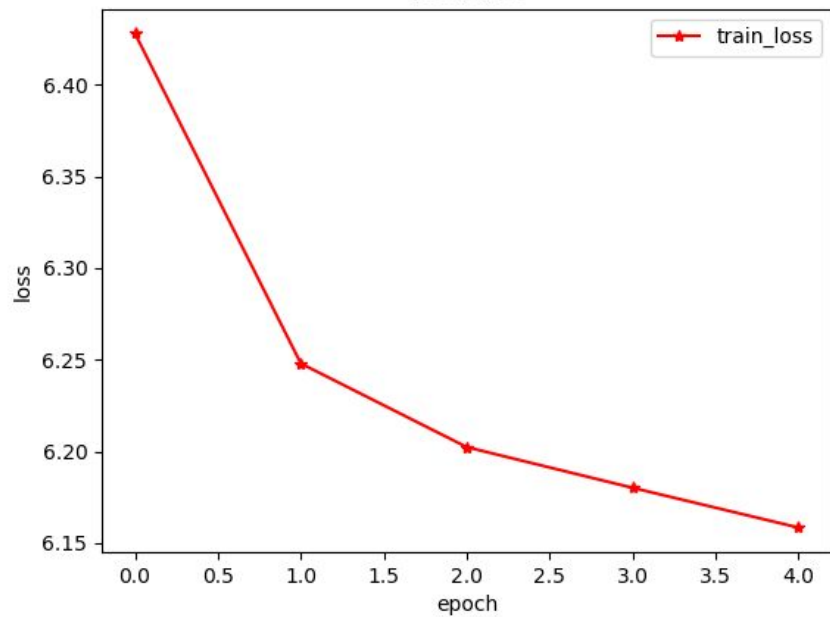
- Single GPU TESLA K40c

# RESULTS

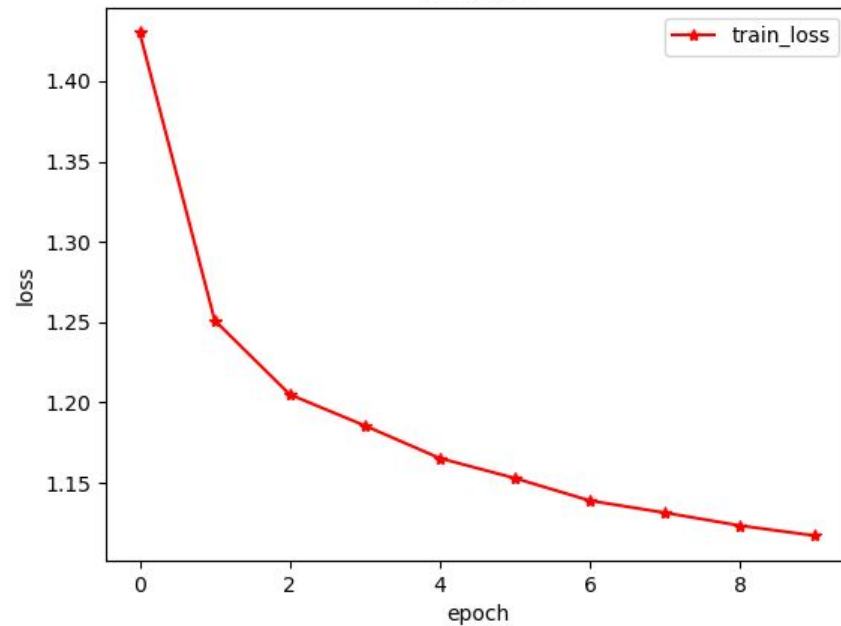| Model | OL | $\mathcal{J} - \mathcal{F}$ mean | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|---|
| PremVOS [13] | Yes | 71,6 | 67,5 | 75,7 |
| MRF [2] | Yes | 67,5 | 64,5 | 70,5 |
| OnAVOS [22] | Yes | 56,5 | 53,4 | 59,6 |
| FeelVOS [21] | No | 57,8 | 55,2 | 60,5 |
| RGMP [16] | No | 52,9 | 51,4 | 54,4 |
| RVOS [20] | No | 50,3 | 47,9 | 52,6 |
| Ours | No | 42.6 | 40.6 | 44.6 |
| Ours w/o AL | No | 44.3 | 42.9 | 45.7 |
| $RVOS_{w/o}$ | No | 33,6 | 32,1 | 35,0 |
| $Ours_{w/o}$ | No | 36,1 | 33,9 | 38,4 |

Table 1. Comparison against state-of-the-art results for semi-supervised video object segmentation in DAVIS-2017 test-dev. OL refers to online learning. $RVOS_{w/o}$ and $Ours_{w/o}$ are the models without the pre-training on YouTube-VOS dataset. Ours w/o AL refers to our model pretrained without the auxiliary loss
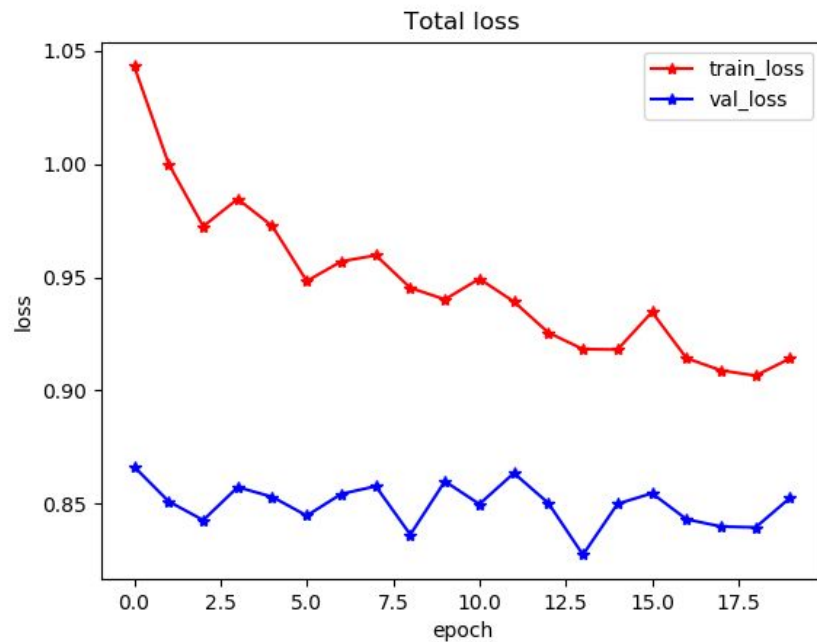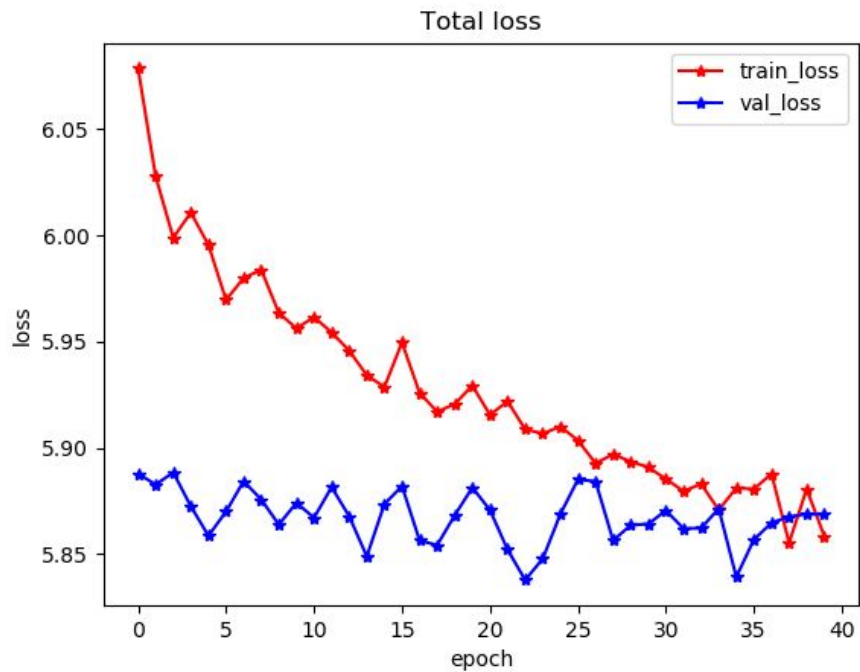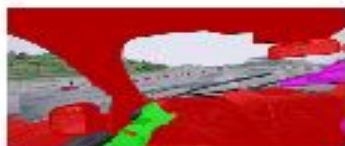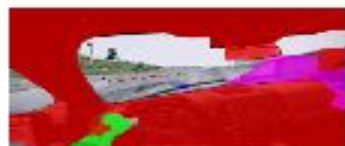
# RESULTS

# RESULTS

# ABLATION STUDY

| Model | FT | $\mathcal{J} - \mathcal{F}$ mean | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|---|
| Ours | No | 41.5 | 40.1 | 42.9 |
| Ours w/o AL | No | 42.9 | 40.8 | 44.9 |
| Ours | Yes | 42.6 | 40.6 | 44.6 |
| Ours w/o AL | Yes | 44.3 | 42.9 | 45.7 |

Table 2. Comparison between model which are just pretrained and those that are fine-tuned in DAVIS dataset.All the results are on DAVIS-2017 test-dev set. FT refers to fine-tuning on DAVIS dataset.Ours w/o AL refers to our model without the auxiliary loss

# ABLATION STUDY

| Model | $r_g$ | $\mathcal{J} - \mathcal{F}$ mean | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|---|
| Ours w/o NLB | 1 | **35,6** | **34,2** | **37,0** |
| Ours w/o NLB | 2 | 31,7 | 30,1 | 33,3 |

Table 3. Ablation study about number of reset gates ($r_g$) on GRU implementation in DAVIS 2017 dataset without pre-training in YouTube-VOS.

# ABLATION STUDY

| Recurrence | NLB | $\mathcal{J} - \mathcal{F}$ mean | $\mathcal{J}$ | $\mathcal{F}$ |
|---|---|---|---|---|
| GRU | Res3 | **35,4** | **33,5** | **37,4** |
| GRU | Res4 | 34,6 | 32,4 | 36,8 |
| GRU | Res5 | 30,4 | 28,9 | 32,0 |
| LSTM | Res3 | 31,7 | 30,0 | 33,4 |
| LSTM | Res4 | 31,2 | 29,8 | 32,5 |
| LSTM | Res5 | 32,1 | 30,6 | 33,7 |

Table 4. Ablation study about the effect of using the NLB module after different parts of the decoder architecture with different kind of recurrence in DAVIS 2017 dataset without pre-training in YouTube-VOS

# ABLATION STUDY

| Model | NLB | $\mathcal{J} - \mathcal{F}$ mean | $\mathcal{J}$ | $\mathcal{F}$ |
|-------|-----|-----------|------|------|
| STESA-GRU | Res3 | **36,1** | **33,9** | **38,4** |
| STESA-GRU$_2$ | Res4 | 33,7 | 31,9 | 35,6 |

Table 5. Ablation study about the effect of using the NLB module with auxiliary loss after different positions in DAVIS 2017 dataset without pre-training in YouTube-VOS

# CONCLUSIONS

- End-to-end trainable spatio temporal recurrent network with a self-attention module

- Comparable results to method which not use online learning in DAVIS test-dev set

- It is necessary to prove training with different learning rate values

- Evaluate the model in YouTube-VOS validation set