1 We thank all reviewers for their constructive and valuable feedback and are delighted to receive an overall positive
2 acceptance of our work. Furthermore, we will integrate all changes according to your suggestions and questions.

3 **Effectiveness (R1, R2)** In the context of our work, efficiency refers to the runtime of the test time optimiza-
4 tion which directly translates to inference time. In order to achieve high segmentation performance previous
5 fine-tuning approaches suffer from unfeasible high runtime due to many fine-tuning epochs (up to 1000). Our
6 approach reduces the number of epochs drastically by meta learning the initialization and learning rates. Fur-
7 thermore, in Figure 1 (which will be included in the supplementary) we illustrate the performance gains and in-
8 vertedly reduction in FPS of our e-OSVOS framework for increasing number of fine-tuning epochs. For many
9 fine-tuning epochs, the actual inference time on the frames is neglectable with respect to the fine-tuning. For
10 longer sequences the amortization of the fine-tuning time is higher, however, these sequences usually face addi-
11 tional challenges due to changing object appearance which can be tackled by the presented online adaptation.
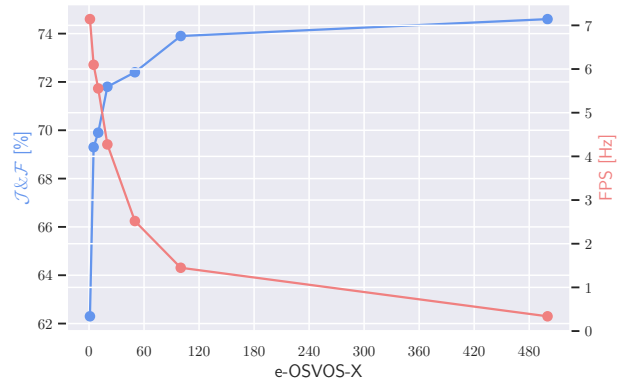12

13 **State-of-the-art (R1, R4)** We understand any potential
14 doubt of our overall impact on the field of VOS as we
15 merely provide a new state-of-the-art for fine-tuning ap-
16 proaches. However, considering publications of the recent
17 years, the VOS community has collectively deemed fine-
18 tuning as unfeasible. Therefore, we hope our demonstra-
19 tion of an efficient fine-tuning approach enabled through
20 meta learning will have a substantial impact and ignite
21 future research based on the release of our code base.

22 **Learning (neuron-level) learning rates (R1)** We pre-
23 dict test time learning rates for fully-connected and convo-
24 lutional layers, each consisting of neurons with a weight
25 vector and a scalar bias. To account for convolutional
26 neurons with a spatial weight tensor (also referred to as
27 kernels), we will rename *weight vector* to *weight tensor*.



Figure 1: Evaluation on the DAVIS 2017 validation set.

28 **(R3)** Learning and deriving a gradient with respect to the learning rates $\lambda$ is analogous to the model initialization $\theta_f^0$.
29 Hence, we formulate the meta optimization of the optimization $g$ for the joint set of its parameters $\theta_g = \{\theta_f^0, \lambda\}$. The
30 gradient flows from the loss (Equation (4)) to each of the $T$ parameter updates (Equation (5)) of the inner fine-tuning
31 loop. The connection between the inner and outer optimization is also illustrated in Algorithm 1 of the supplementary.
32 As the SGD update consists only of differentiable operations, gradients with respect to the learning rates can be derived
33 analogous to the derivation for the initial parameters. It should be noted, that these gradients are with respect to a
34 different loss ($\mathcal{L}_{seg}(\mathcal{D}_{test}, \theta_f^T)$) as the ones of the inner gradient ($\mathcal{L}_{seg}(\mathcal{D}_{train}, \theta_f^t)$). It is a common approach to meta
35 learn the initialization and learning rate(s) jointly. We refer the reader to [2, 1] for further insights. **(R2)** At test time,
36 the set of learning rates does not change and is the same for all sequences. However, it is interesting to observe which
37 neurons are updated with particularly small, e.g., biases of last layers, or large, e.g., FC6 of the box head, learning
38 rates. The FC6 layer of the box head prepares the spatial bounding box features for the regression and classification
39 heads and benefits from a strong adaption to each individual given object. To further illustrate e-OSVOS, we will add a
40 summarized visualization of the overall more than 20000 neuron-level learning rates to the supplementary.

41 **Other comments/suggestions (R1)** In the ablation study (first row of Table 1), we present a Mask R-CNN baseline
42 without meta learning which was pre-trained on ImageNet, COCO segmentation, YouTube-VOS and DAVIS 2017.
43 However, this baseline is not representative for state-of-the-art fine-tuning approaches as we omitted any additional
44 handcrafted test time improvements. For a state-of-the-art fine-tuning approach without meta learning but with bells
45 and whistles, e.g., online adaption, we compare to OnAVOS [3] in Table 3. **(R2)** Upon acceptance we will publish
46 our results on the official DAVIS challenge webpage which provides a tool for visual comparison per sequence, e.g.,
47 blackswan. **(R3)** The *mitigate* in line 7 refers to *shortcomings* and we will improve the understandability of the abstract.
48 The fine-tuning epochs in Table 1 refer to a single update with one image. However, to improve the generalization over
49 a sequence, we train on batches of random transformations of that image. The superscript of $\theta_f$ always implies how
50 the parameters were optimized, e.g., for $T$ update steps or by optimization $g$. In our formulation, the optimization $g$
51 describes a model initialization, a set of learning rates and number of steps $T$. We will clarify this overloading of the
52 superscript in the final version. Furthermore, we observed overfitting without the YouTube-VOS dataset and therefore
53 train on a combination of all three datasets.

54 [1] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your MAML. In *International Conference on Learning Representations*, 2019.

55 [2] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conf. on Machine Learning*, ICML'17, pages 1126–1135.
56 JMLR.org, 2017.

57 [3] Paul Voigtlaender and Bastian Leibe. Online adaptation of convolutional neural networks for video object segmentation. In *BMVC*, 2017.