

Performance of "meta", "meta-act" and "act_sb" models when optimizing multilayer perceptrons (MLPs: 1 hidden layer with 20 units using sigmoid activation function) on MNIST (10 classes).

All models consist of an RNN network (LSTM cells) with 2 layers and a hidden state size of 40 units.

The models were evaluated on 50 and 200 newly sampled MLPs and unrolled for **100** and **200** time steps.

Please note that during training the meta_act and act_sb models were unrolled for a maximum of 100 time steps in case the model did not "cross" the threshold (fixed for meta_act and uniformly sampled for act_sb) for a particular optimizee in order to prevent the model from taking too many steps.

More specific the following models were evaluated:

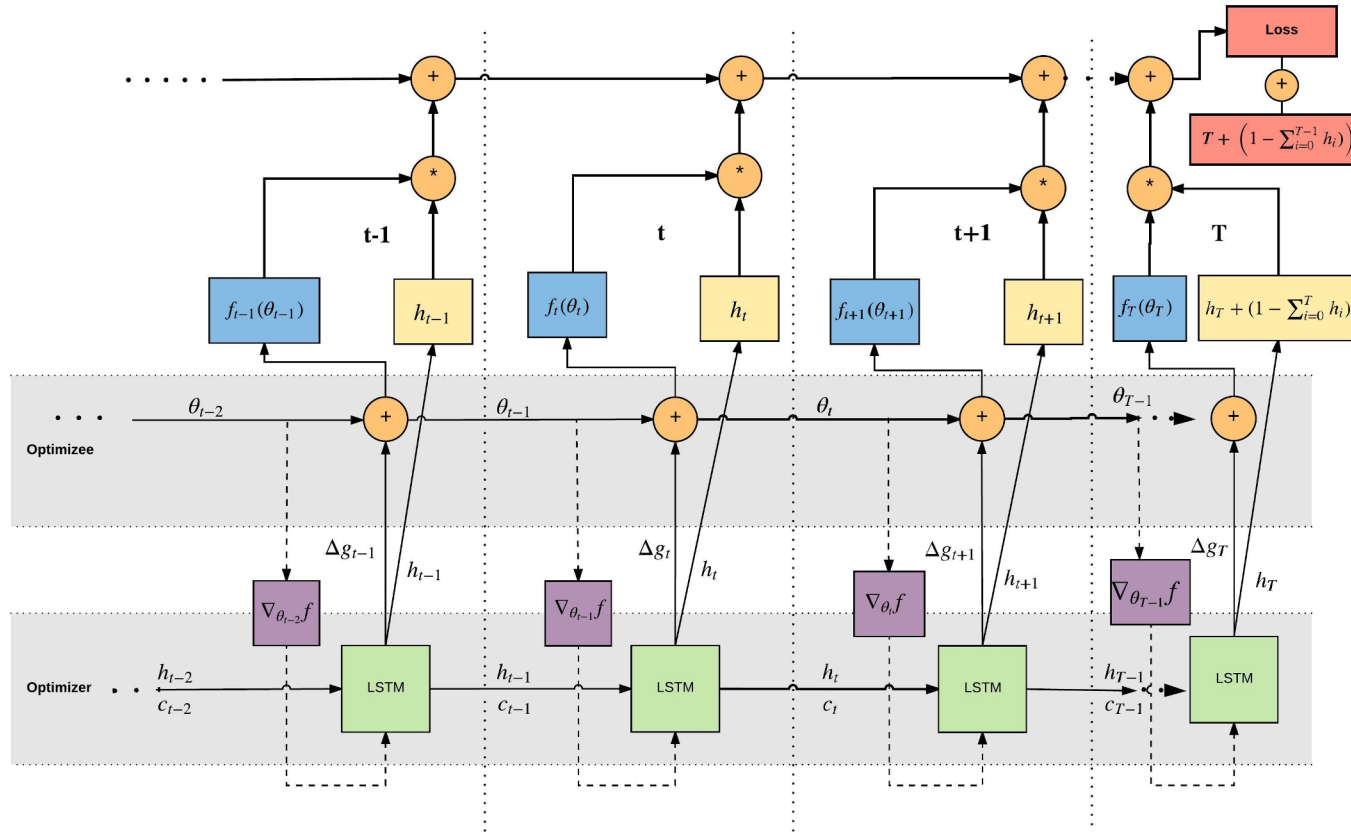
(1) **metaV1**: baseline model from L2L paper where each MLP is trained for a horizon of 50 and 100 time steps ($lr=1e-5$) using truncated BPTT (each 20 steps). Model was once trained for 13 epochs using 78 MLPs/epoch and once for 90 epochs using 16 MLPs per epoch.

(2) **meta_actV1**: combination of meta learner from L2L paper and ACT model from Alex Graves. The model is trained ($lr=5e-5$) with different settings for the hyperparameter tau that influences the trade off between accuracy and speed (it basically scales the penalty term in the loss function). During training the model could not unroll for more than 100 time steps (also experimented with 50 time steps as fixed horizon H). Models were trained for 150 epochs with 16 MLPs/epoch. A fixed threshold of 0.99 is used in order to determine the halting step.

(3) **act_sbV3.2**: difference with meta_actV1 (a) probabilities are calculated based on stick-breaking approach; (b) penalty term consists of a scaled KL-term (geometric prior with shape parameter nu); (c) trained with $lr=1e-5$ for 125 epochs with 16 MLP/epoch; (d) uses a stochastic threshold to sample the halting step (uniform between 0 and 1).

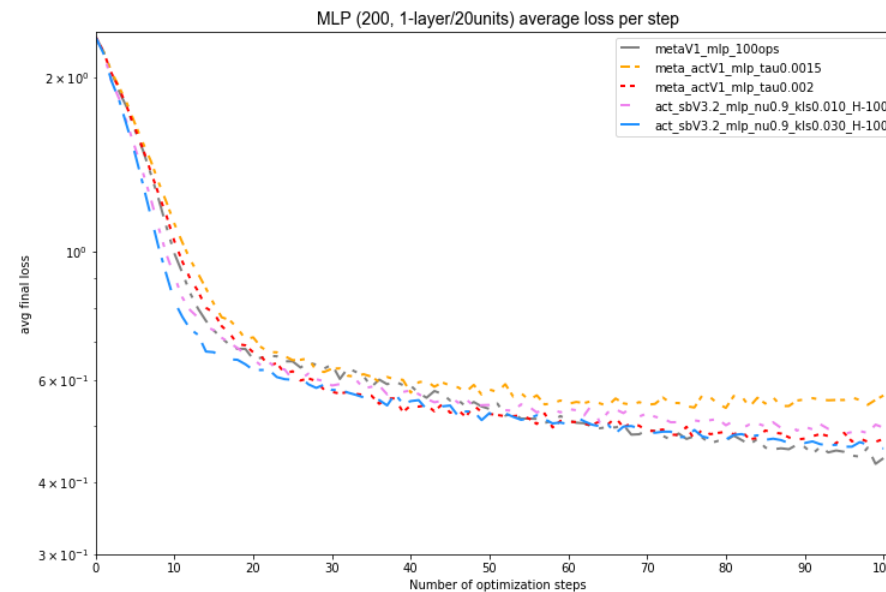
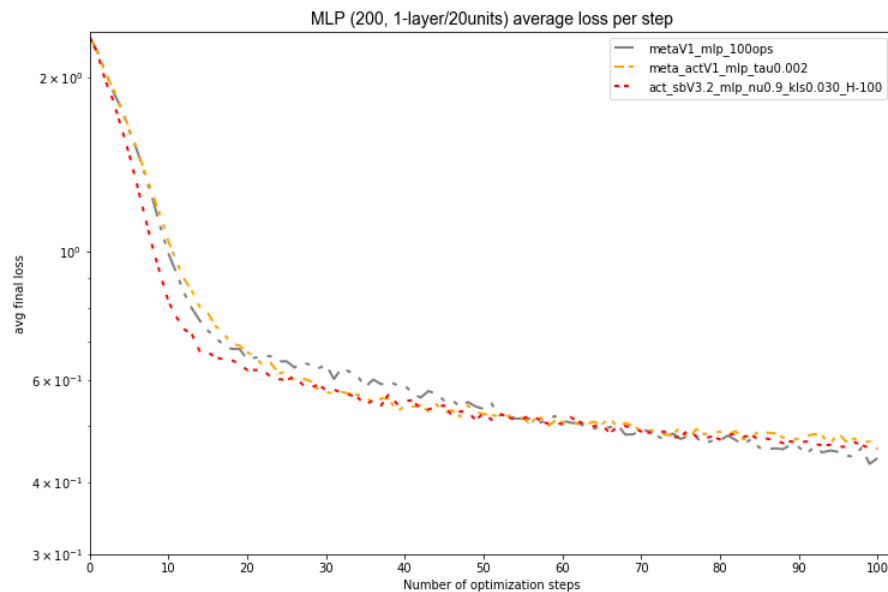
Models were trained with different KL-rescaling: 0.01 and 0.03.

Computational graph meta_act model



Comparison of test performance of metaV1 - meta_actV1 - act_sbV3.2 on 200, 1-layer MLPs

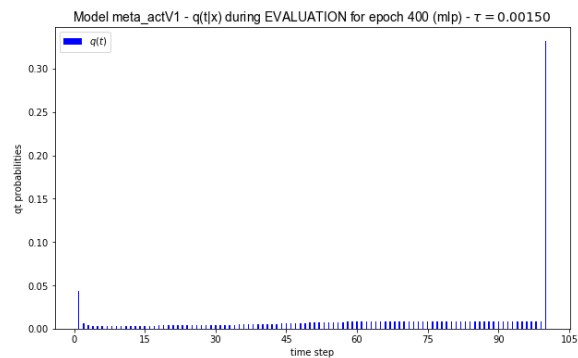
2



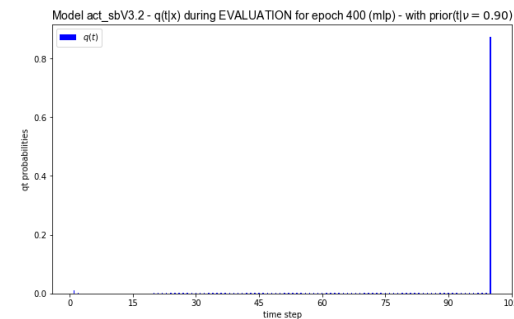
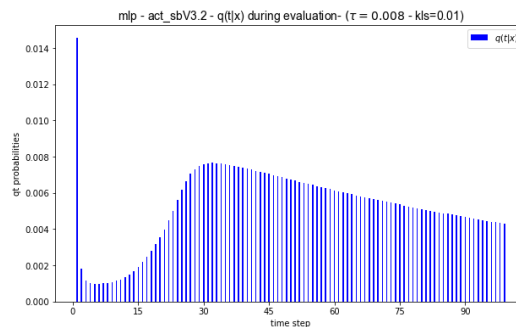
meta_actV1

$q(t|x)$ distributions

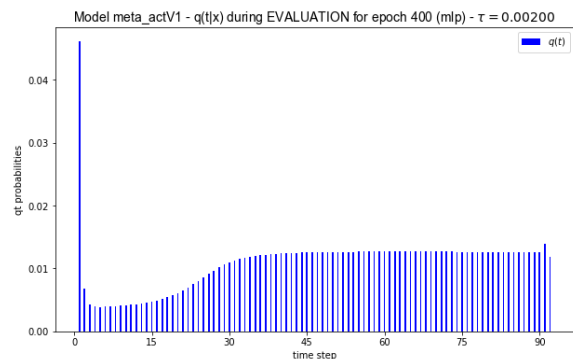
act_sbV3.2



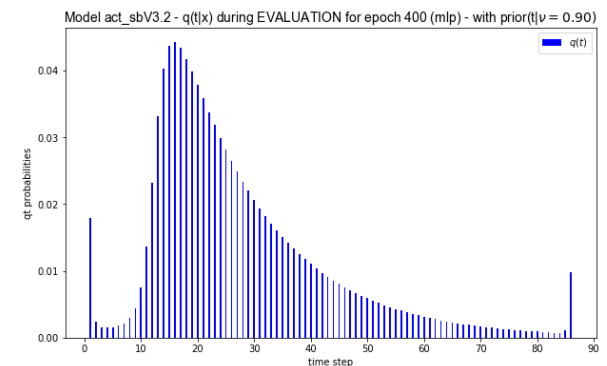
$\tau = 0.0015$



$kls=0.01$



$\tau = 0.002$



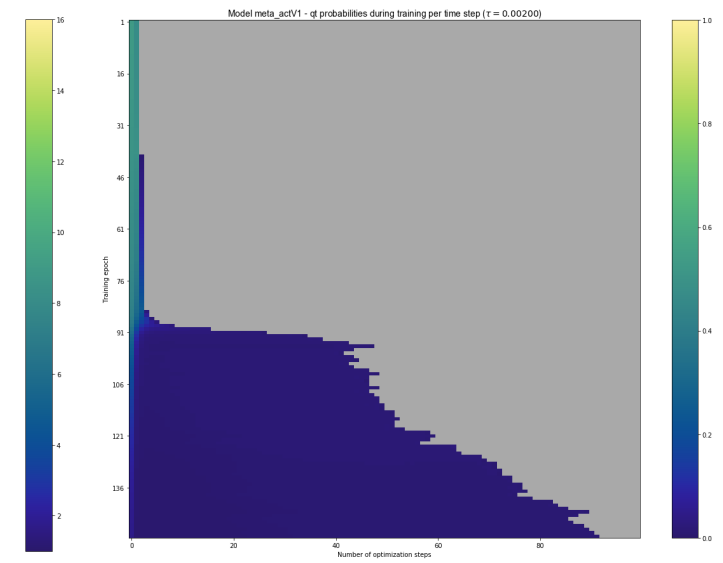
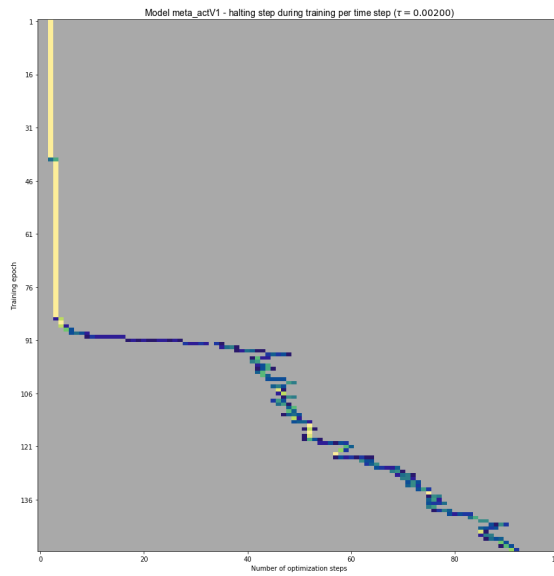
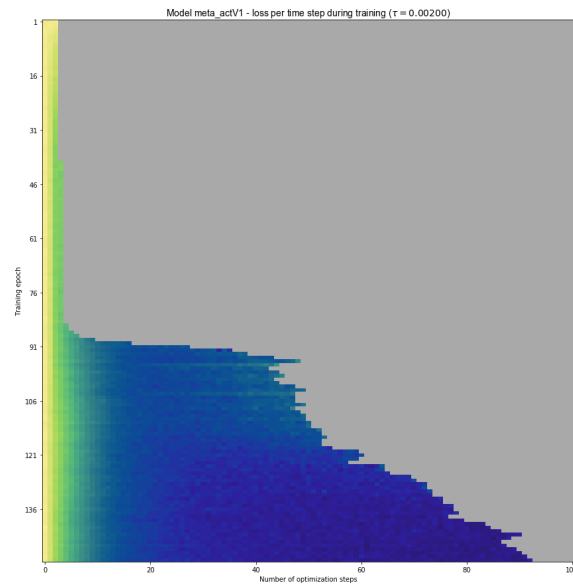
$kls=0.03$

$\tau = 0.002$

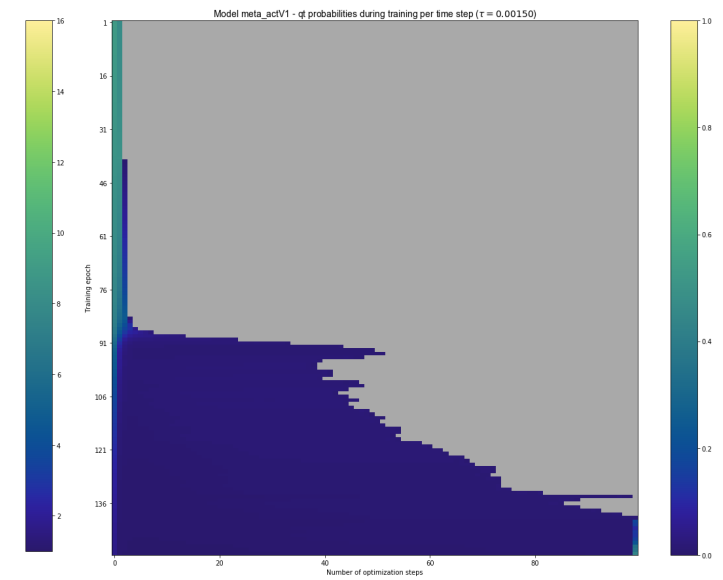
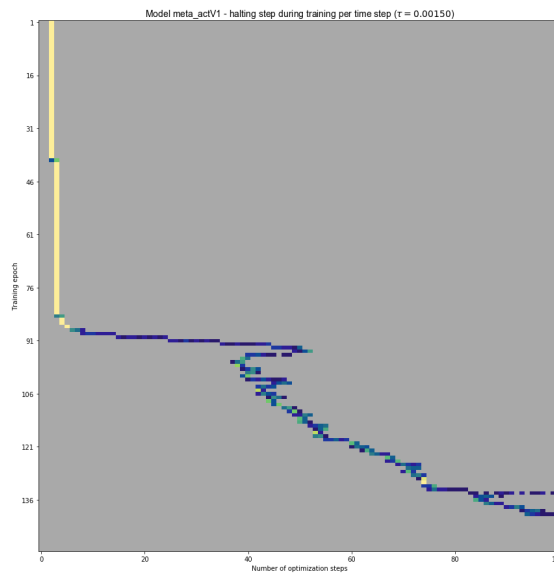
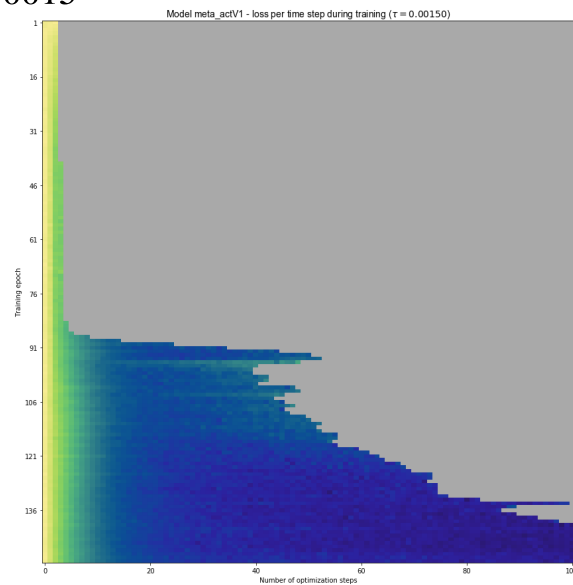
loss

halting step distribution

halting probabilities

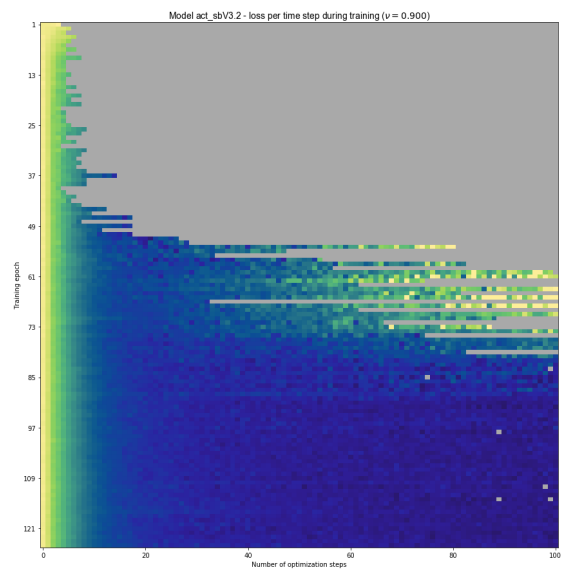


$\tau = 0.0015$

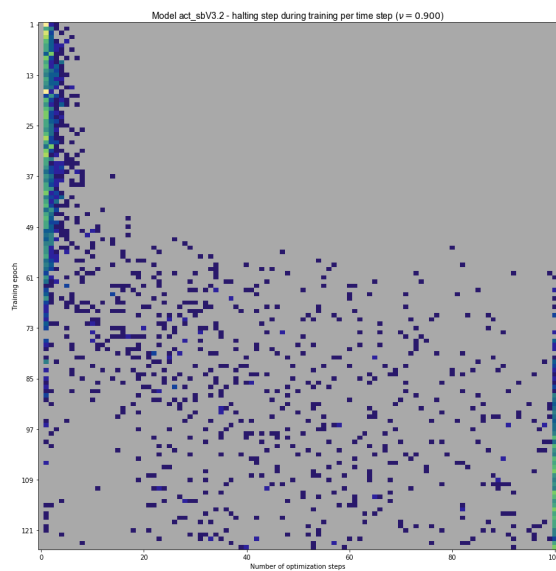


kl-rescaling=0.01

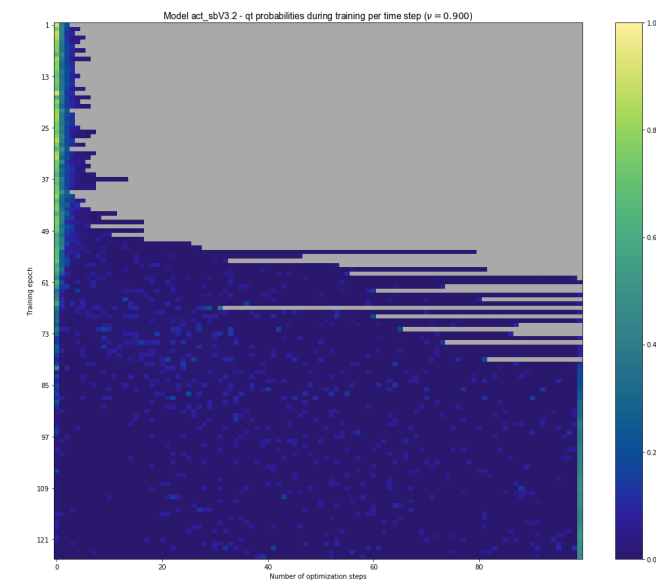
loss



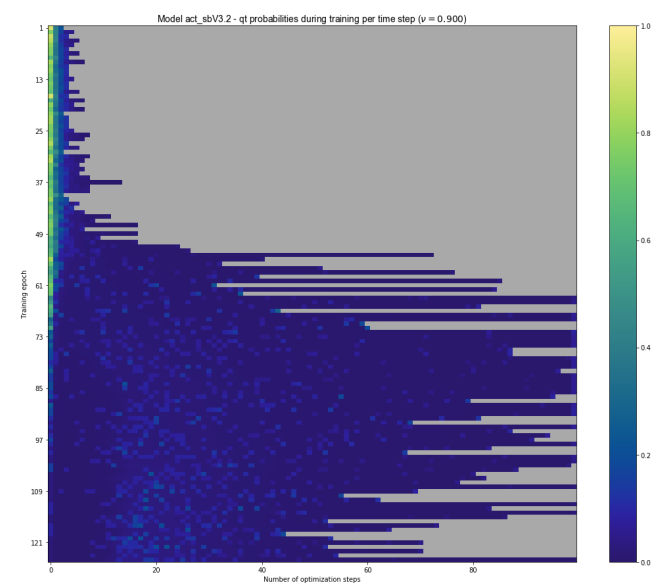
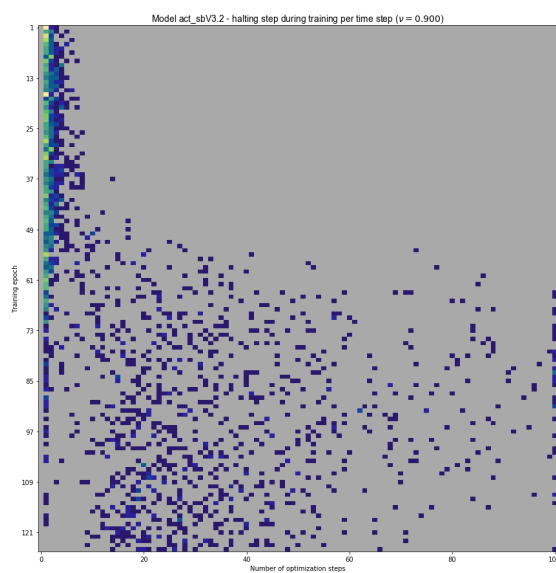
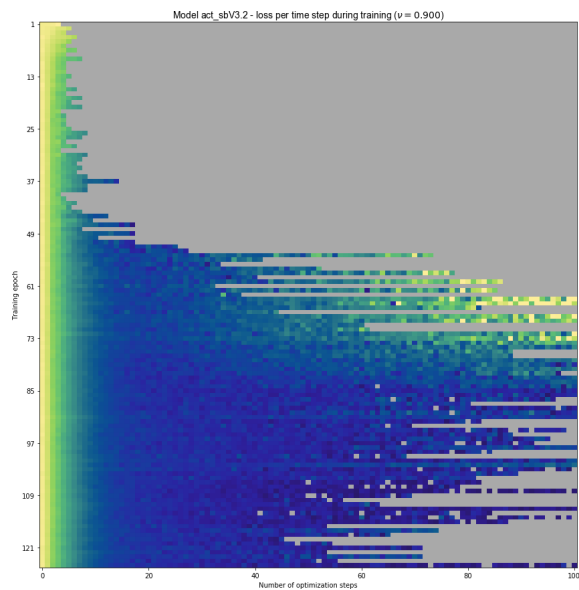
halting step distribution

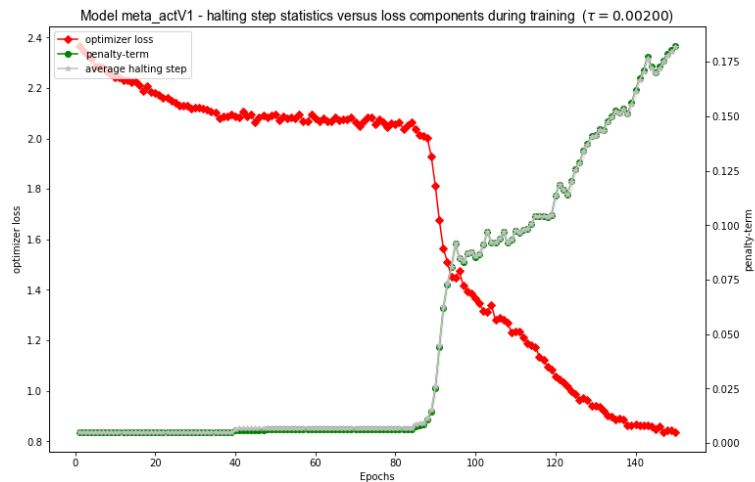


halting probabilities

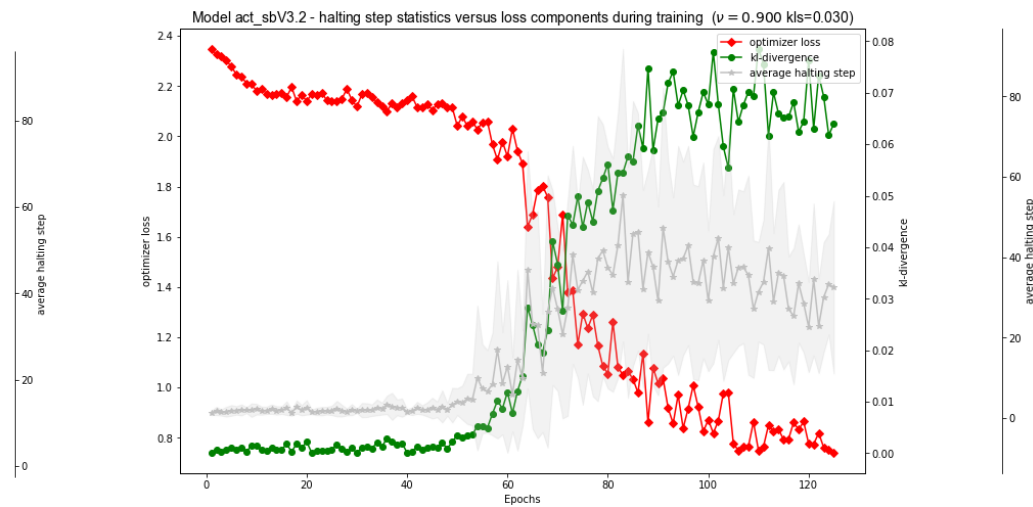


kl-rescaling=0.03

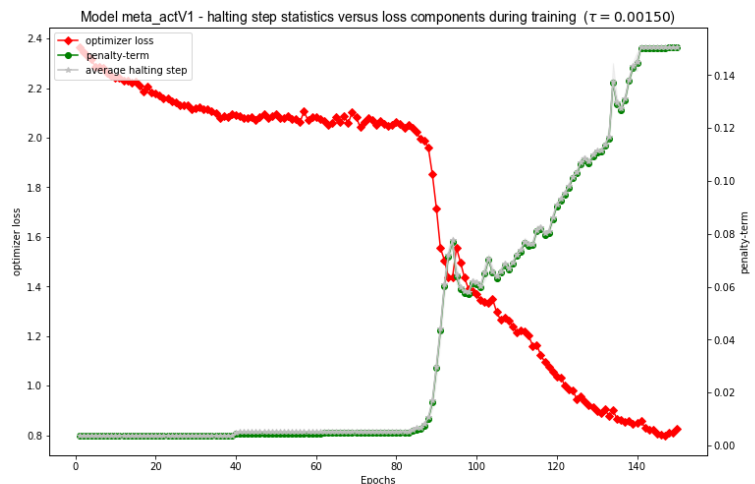




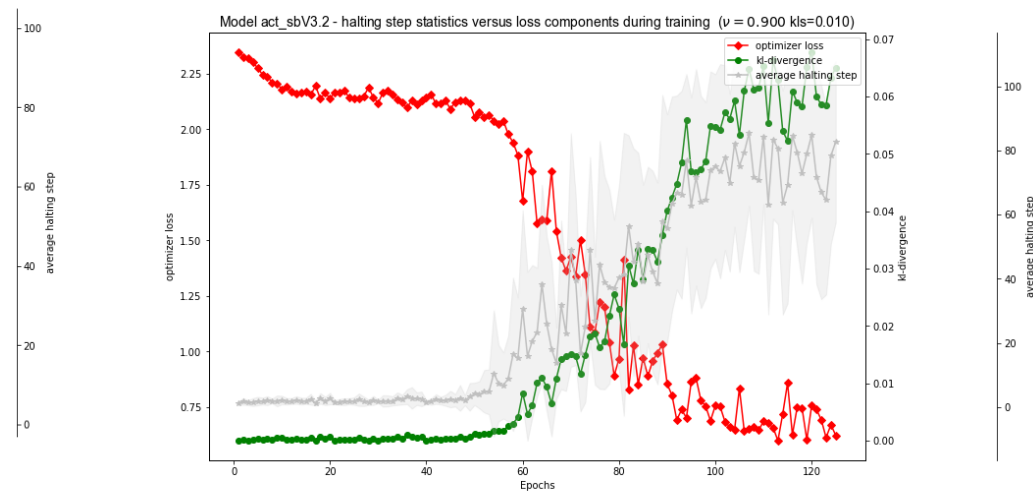
Halting step statistics final epoch: Range(91, 92) / mean=91.2 / stddev=0.4 / median=91.0



Halting step statistics final epoch: Range(14, 100) / mean=32.5 / stddev=21.4 / median=24.0



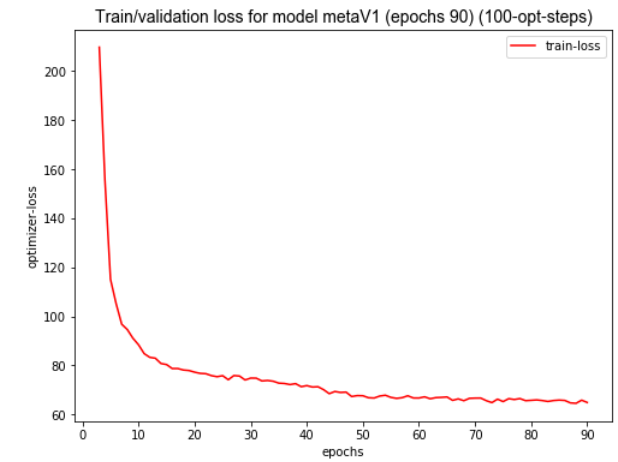
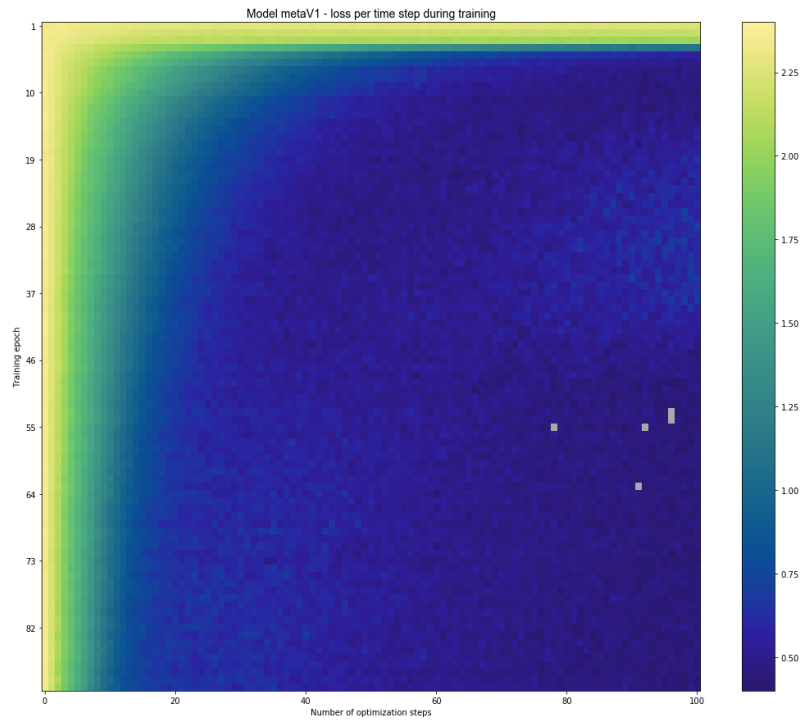
Halting step statistics final epoch: Range(100, 100) / mean=100.0 / stddev=0.0 / median=100.0



Halting step statistics final epoch: Range(39, 100) / mean=82.7 / stddev=25.1 / median=93.0

Learning curves of **metaV1** model

6

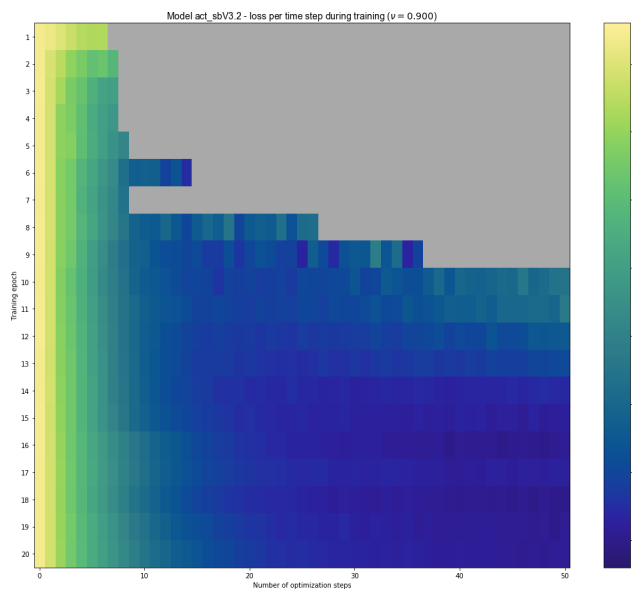


Effect of finite horizon (H) on learning - act_sbV3.2 (kls=0.01) - used **H=50** and **H=100**

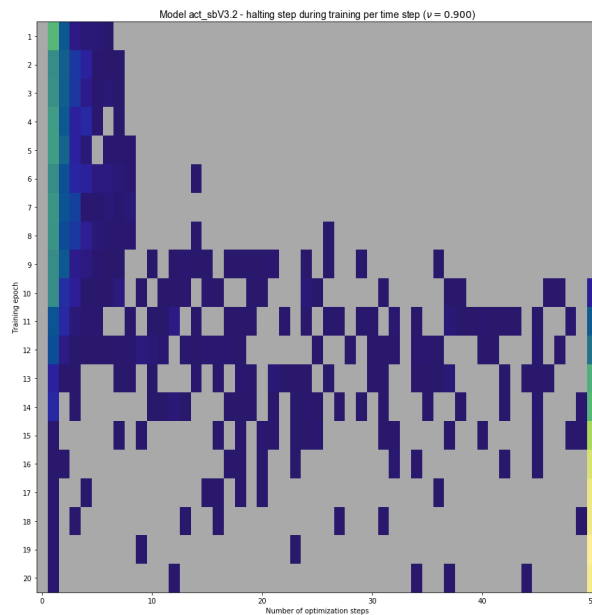
7

With **H = 50** (see page 4 w.r.t. same model trained with H=100)

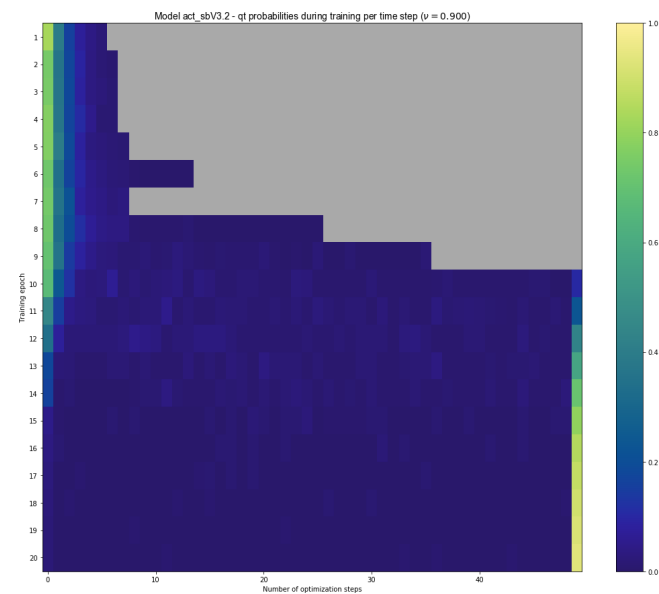
loss



halting step distribution



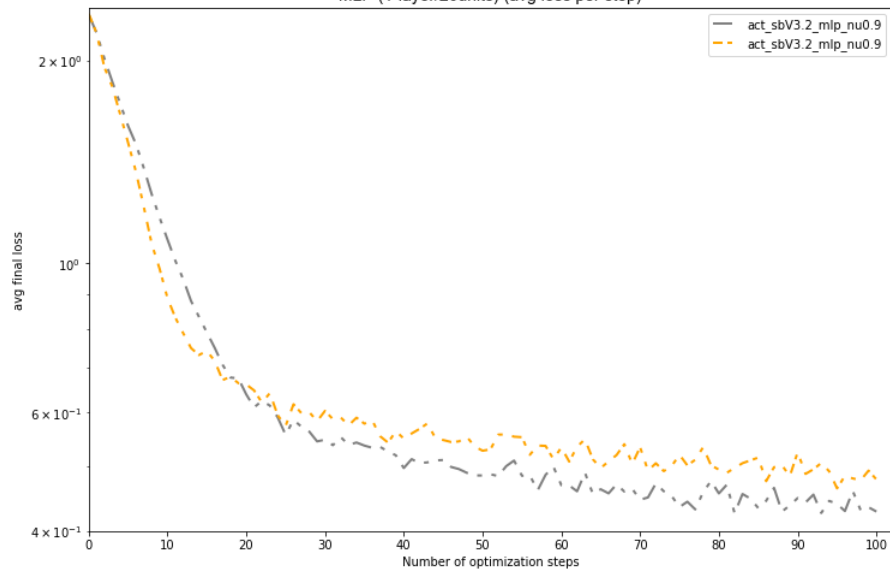
halting probabilities



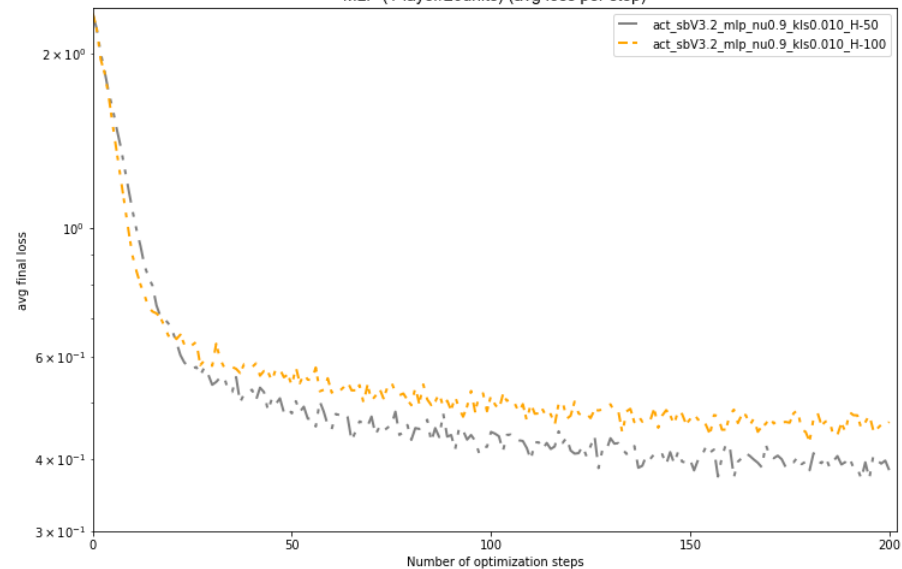
"Yellow" model is trained with **H = 100**

Both models were trained for 2000 MLPs (20x100 and 150x16)

MLP (1-layer/20units) (avg loss per step)

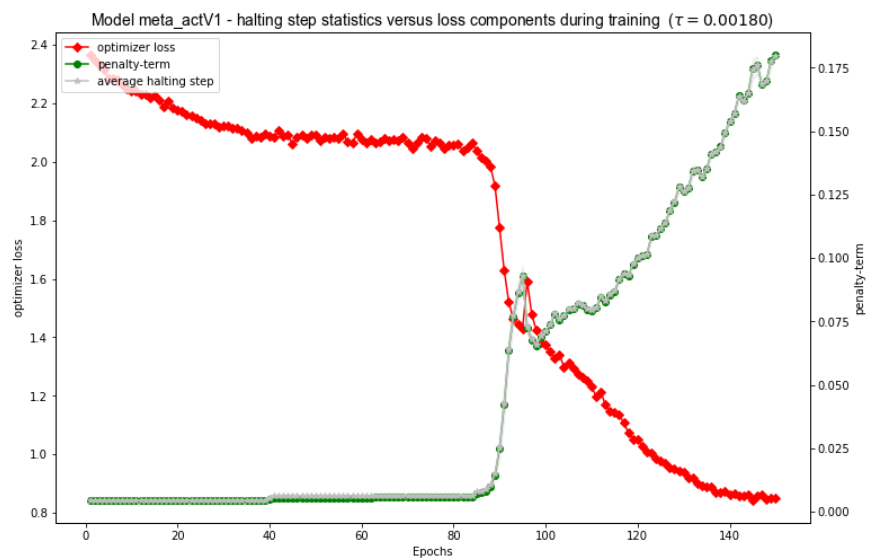
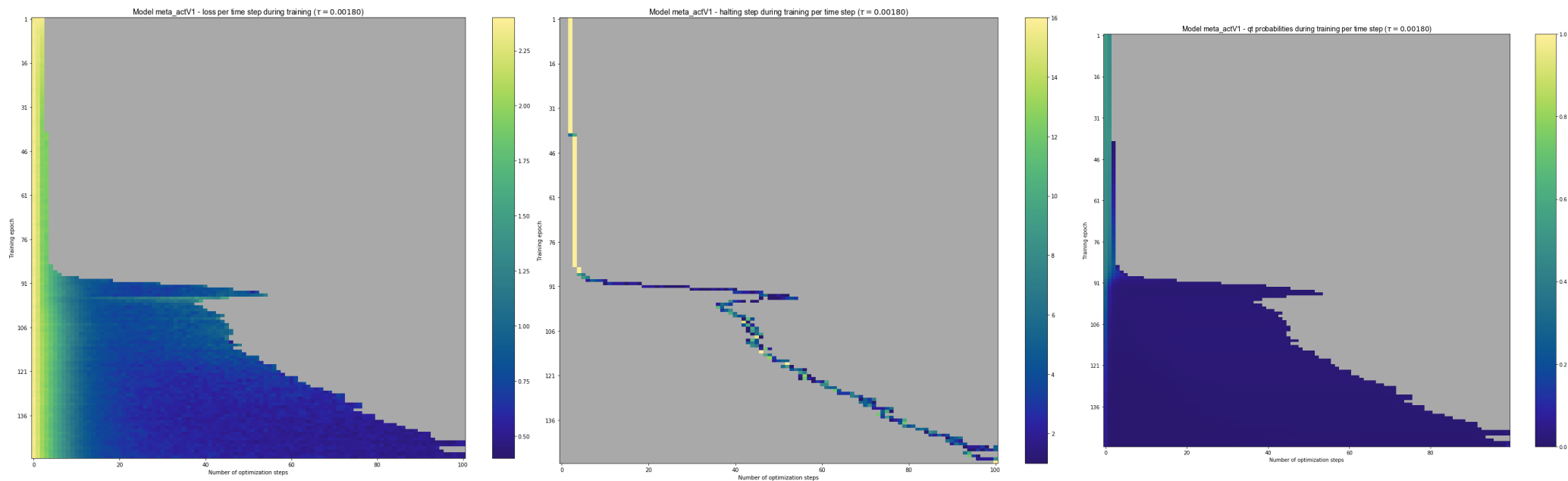


MLP (1-layer/20units) (avg loss per step)



Details training results **meta_actV1** ($\tau=0.018$) to illustrate effect of τ (compare with page 3)

8



Halting step statistics final epoch: Range(100, 100) / mean=100.0 / stddev=0.0 / median=100.0

