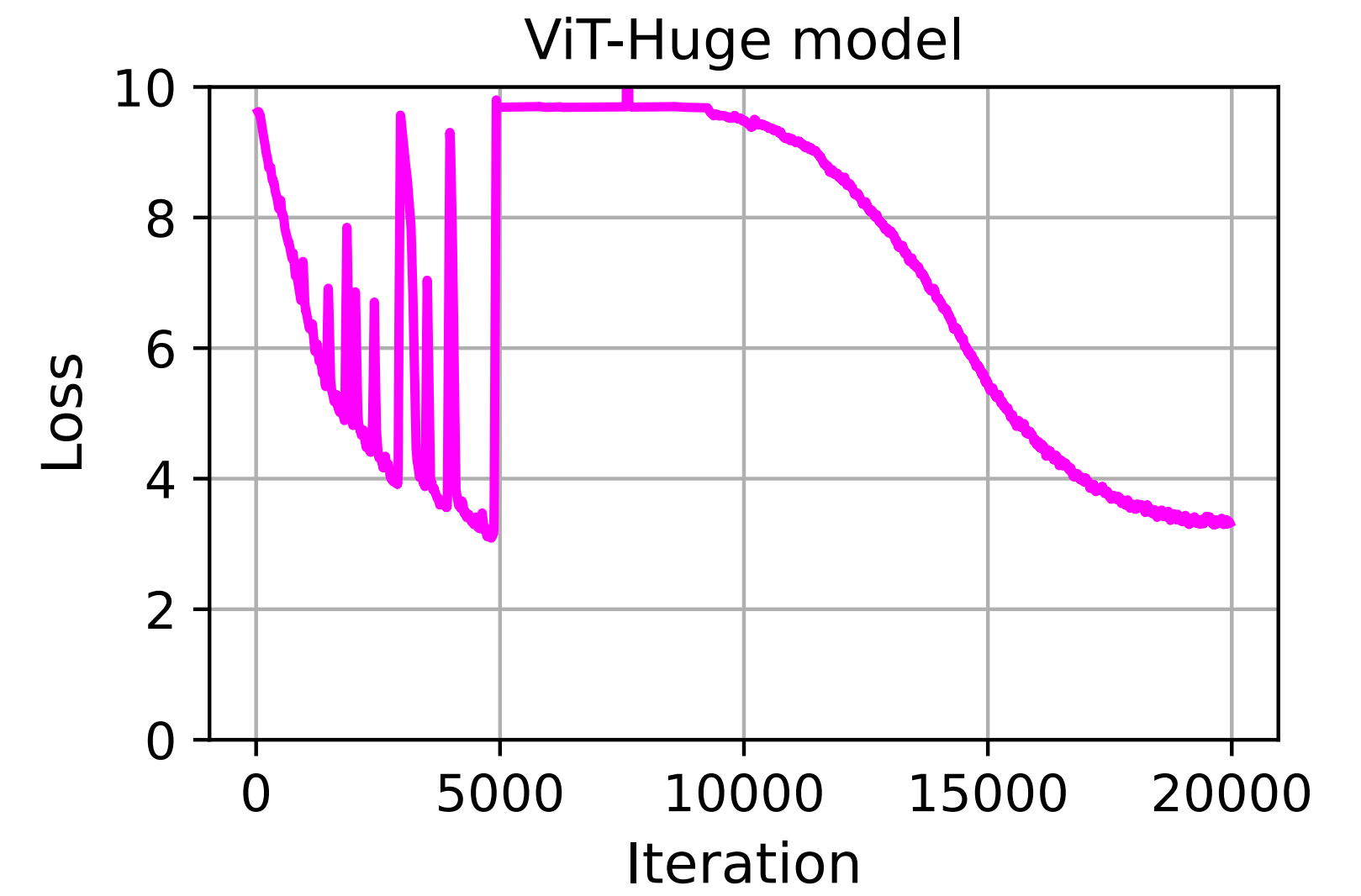
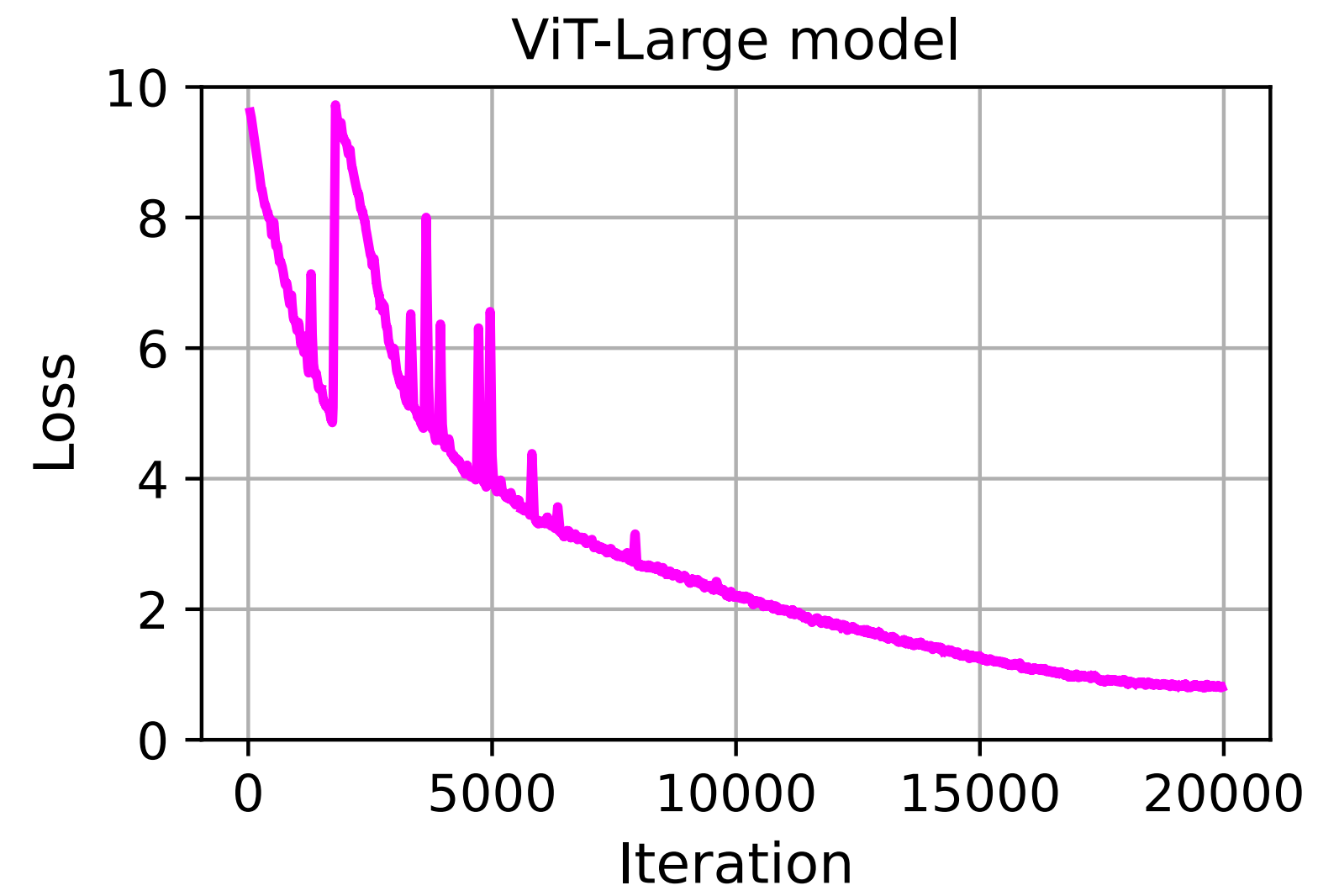
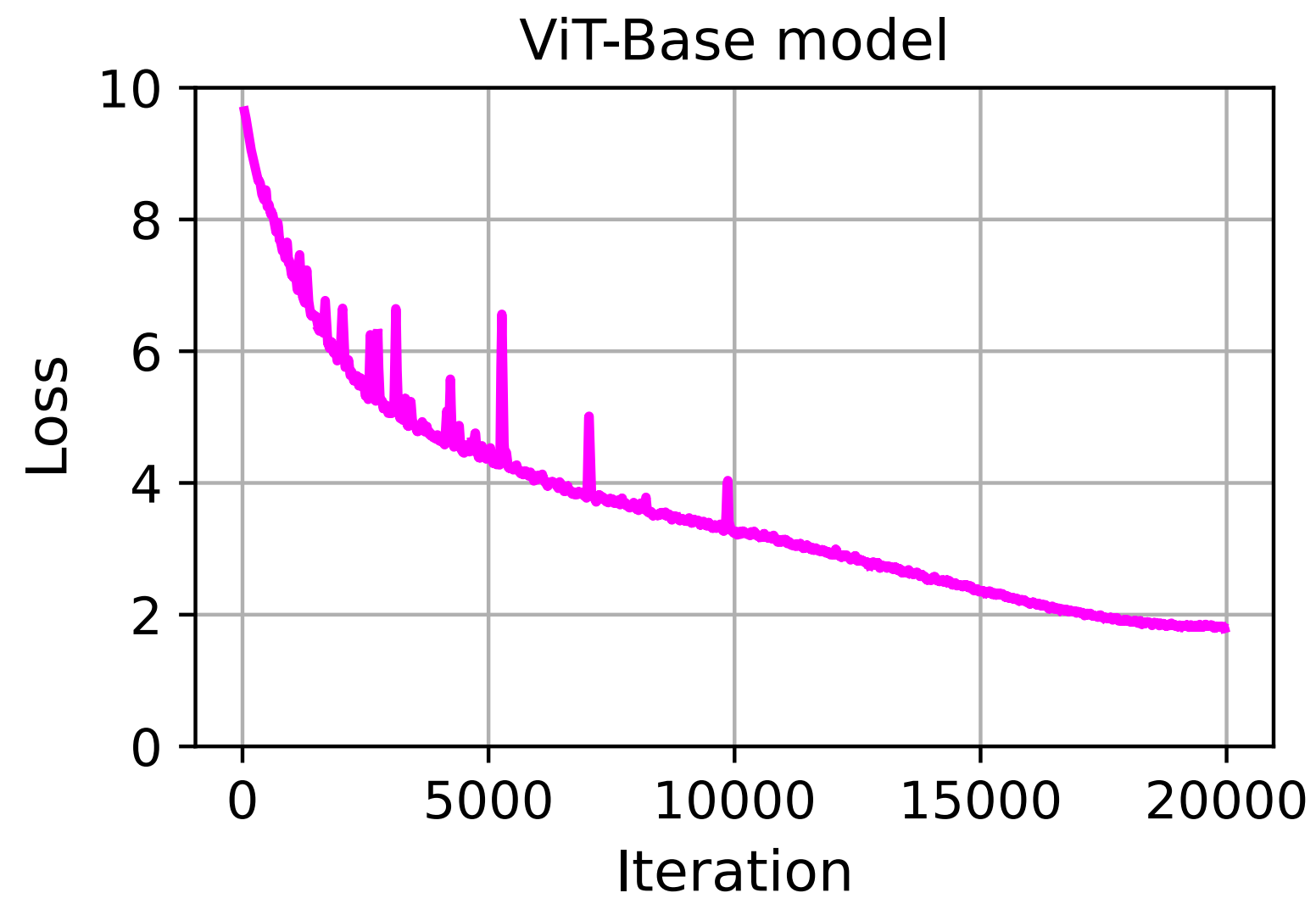


# Adam - what can go wrong and how to fix it

Mitchell

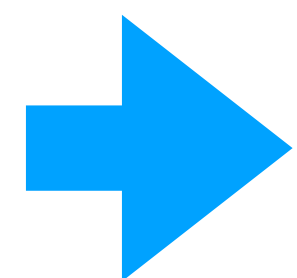
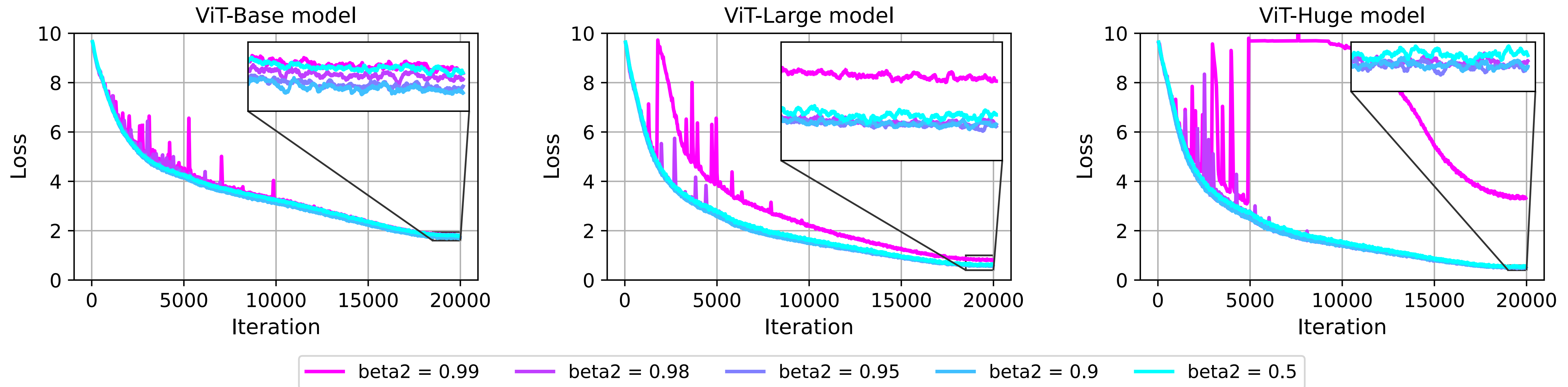
Based on part of an upcoming research project with Tim Dettmers, Luke Zettlemoyer, Ari Morcos, Ali Farhadi, Ludwig Schmidt

# What can go wrong with Adam?



Beta1 = 0.9, Beta2 = 0.99

## Tuning more AdamW beta2 values



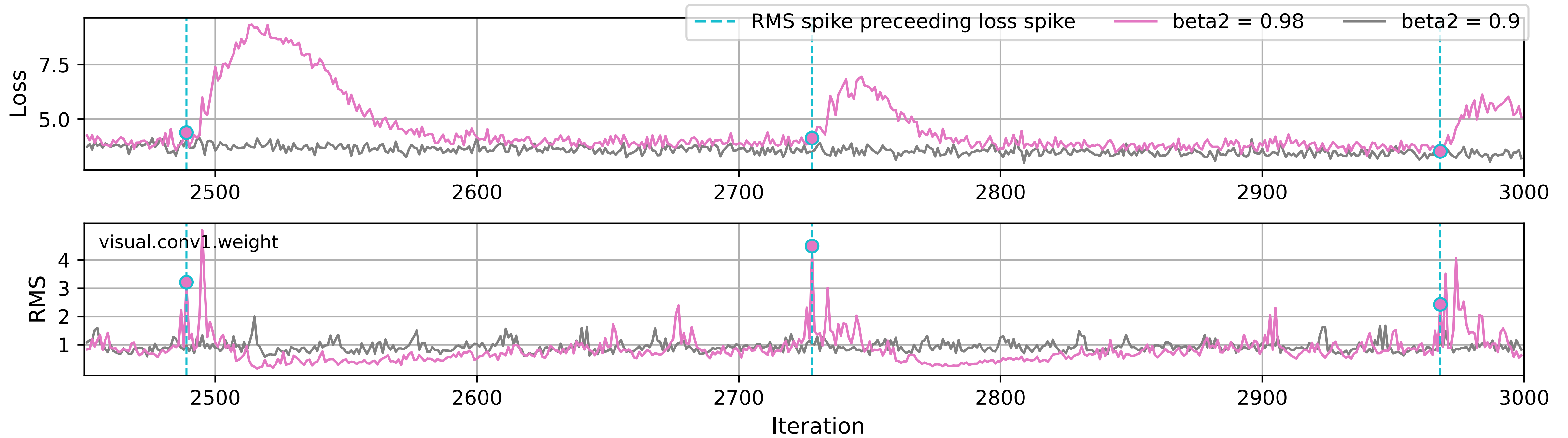
**The second moment estimator becomes out-of-date, leading to updates which are too large.**

Let's track how good of an estimator we have for the squared gradient... the following aggregate quantity should be approx 1

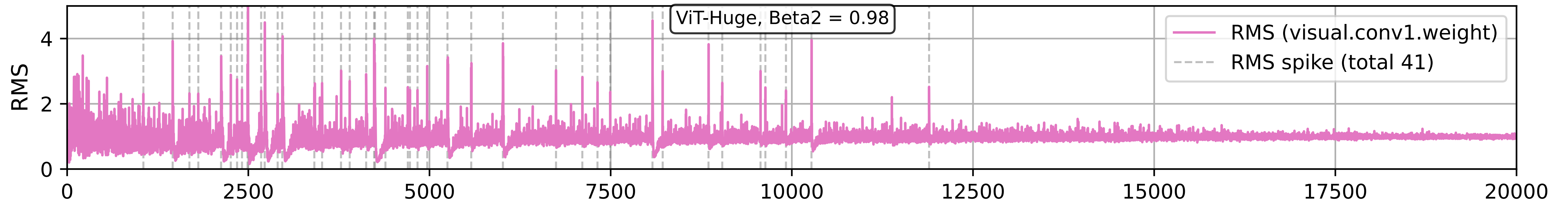
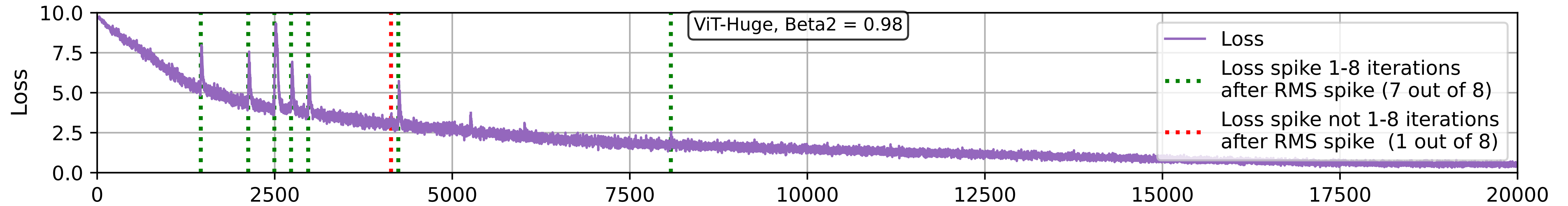
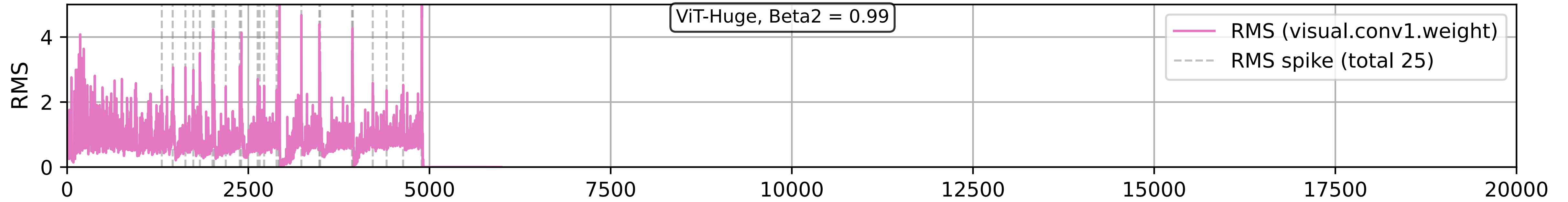
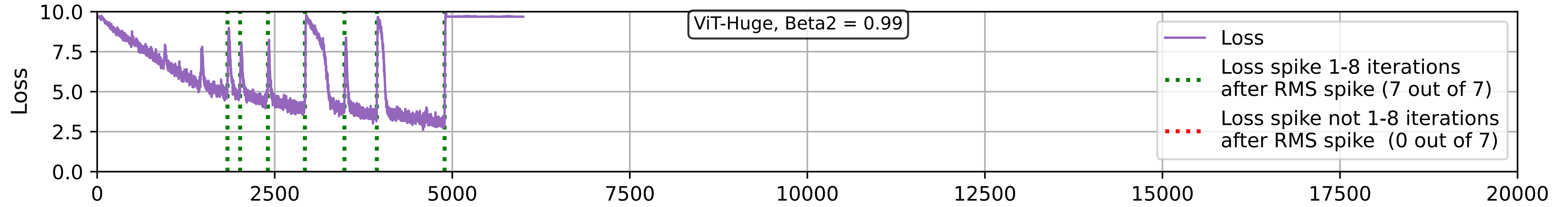
$$\text{RMS}_t = \sqrt{\mathbb{E} \left[ \frac{g_t^2}{u_t} \right]}$$

Can we search for a causal relationship between an RMS spike and a loss spike?

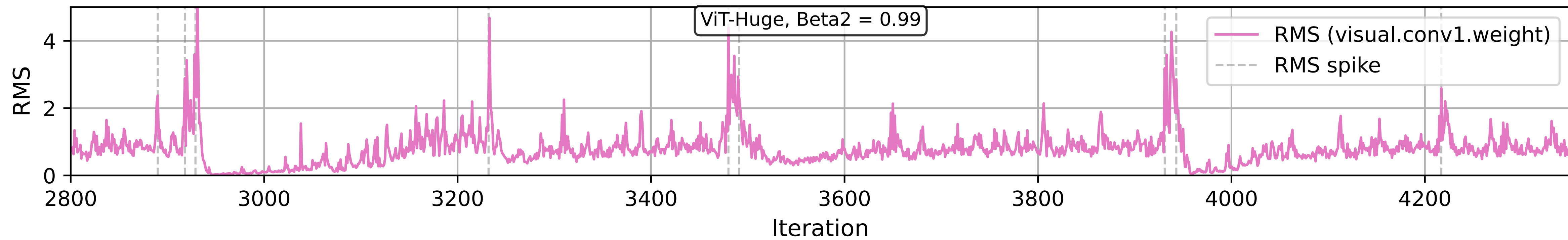
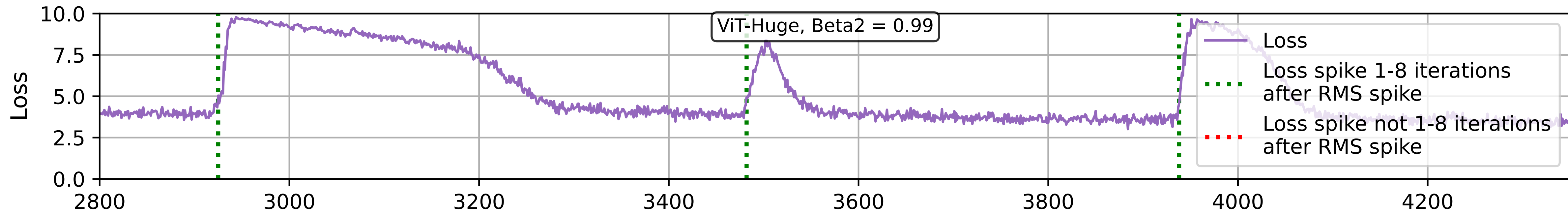
## RMS spikes in the patch embedding layer come before loss spikes



**This is just one example... let's come up with heuristic for documenting loss/rms spikes and use this for analysis**

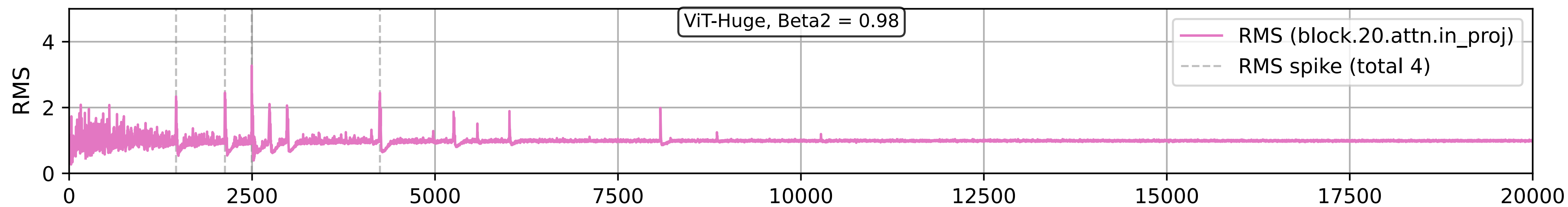
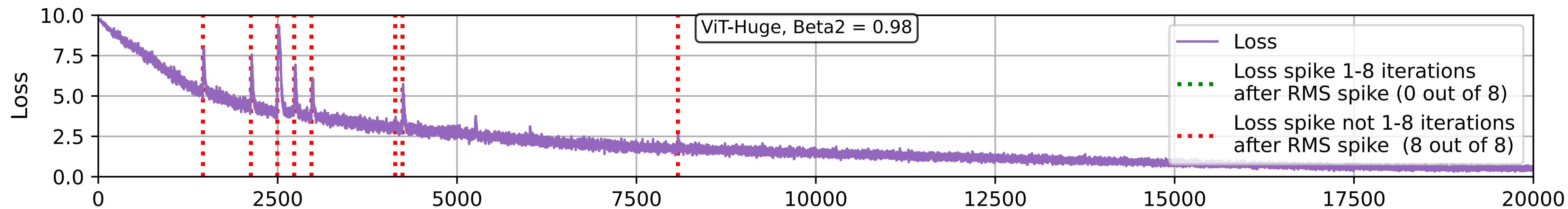
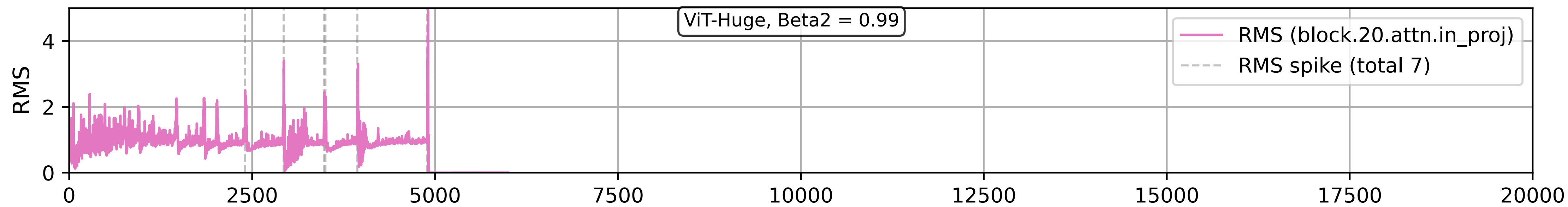
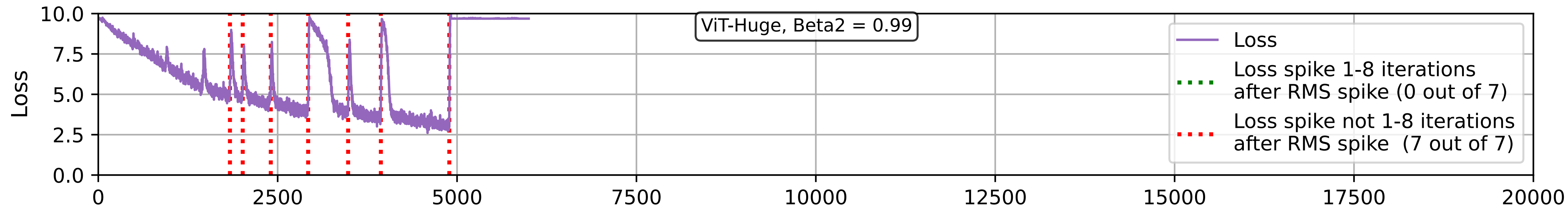


**Chance of a loss spike randomly following an RMS spike is 1.0%**





**For later layers, RMS spikes are postdictive not predictive of loss spikes**



Iteration

- Has this ever been encountered before?
  - Yes, Shazeer and Stern found an out of date second moment estimator when developing AdaFactor and running experiments without warm-up
  - They propose a fix called “update clipping”
- Does anyone use “update clipping”?
  - No, because people do not get good performance with AdaFactor, but this is for other reasons (most likely the factored moments)
- So... let’s try AdamW + update clipping and refer to the result as StableAdamW

## Slow down learning (reduce learning rate) if second moment estimator out of date

---

**Algorithm 2** StableAdamW ( $\{\alpha_t\}_{t=0}^T, \beta_1, \beta_2, \epsilon$ )

---

$$v_0, u_0 = \mathbf{0}$$

**for**  $t = 1$  **to**  $T$  **do**

$$g_t = \nabla f(\theta_t)$$

// apply correction term to debias moving avg.

$$\hat{\beta}_1 = \beta_1 \cdot \frac{1 - \beta_1^{t-1}}{1 - \beta_1^t}$$

$$\hat{\beta}_2 = \beta_2 \cdot \frac{1 - \beta_2^{t-1}}{1 - \beta_2^t}$$

// update moving averages

$$v_t = \hat{\beta}_1 v_{t-1} + (1 - \hat{\beta}_1) g_t$$

$$u_t = \hat{\beta}_2 u_{t-1} + (1 - \hat{\beta}_2) g_t^2$$

// for implementation convenience, the steps

// below occur independently for each tensor

$$\mathbf{RMS}_t = \sqrt{\mathbb{E}[g_t^2 / u_t]}$$

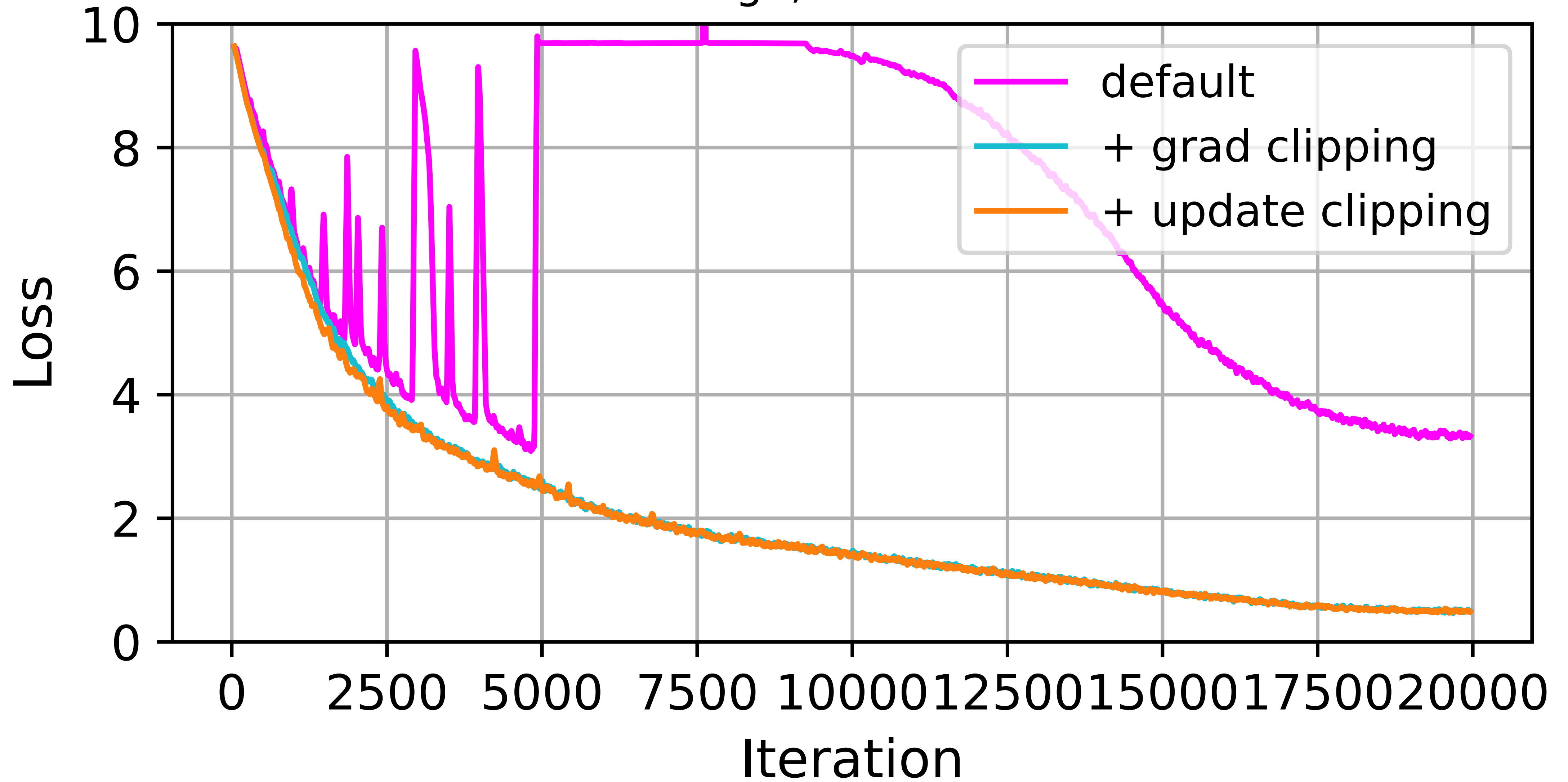
// update parameters

$$\eta_t = \alpha_t / \max(1, \mathbf{RMS}_t)$$

$$\theta_t = \theta_{t-1} - \alpha_t \lambda \theta_{t-1} + \eta_t v_t / (\sqrt{u_t} + \epsilon)$$

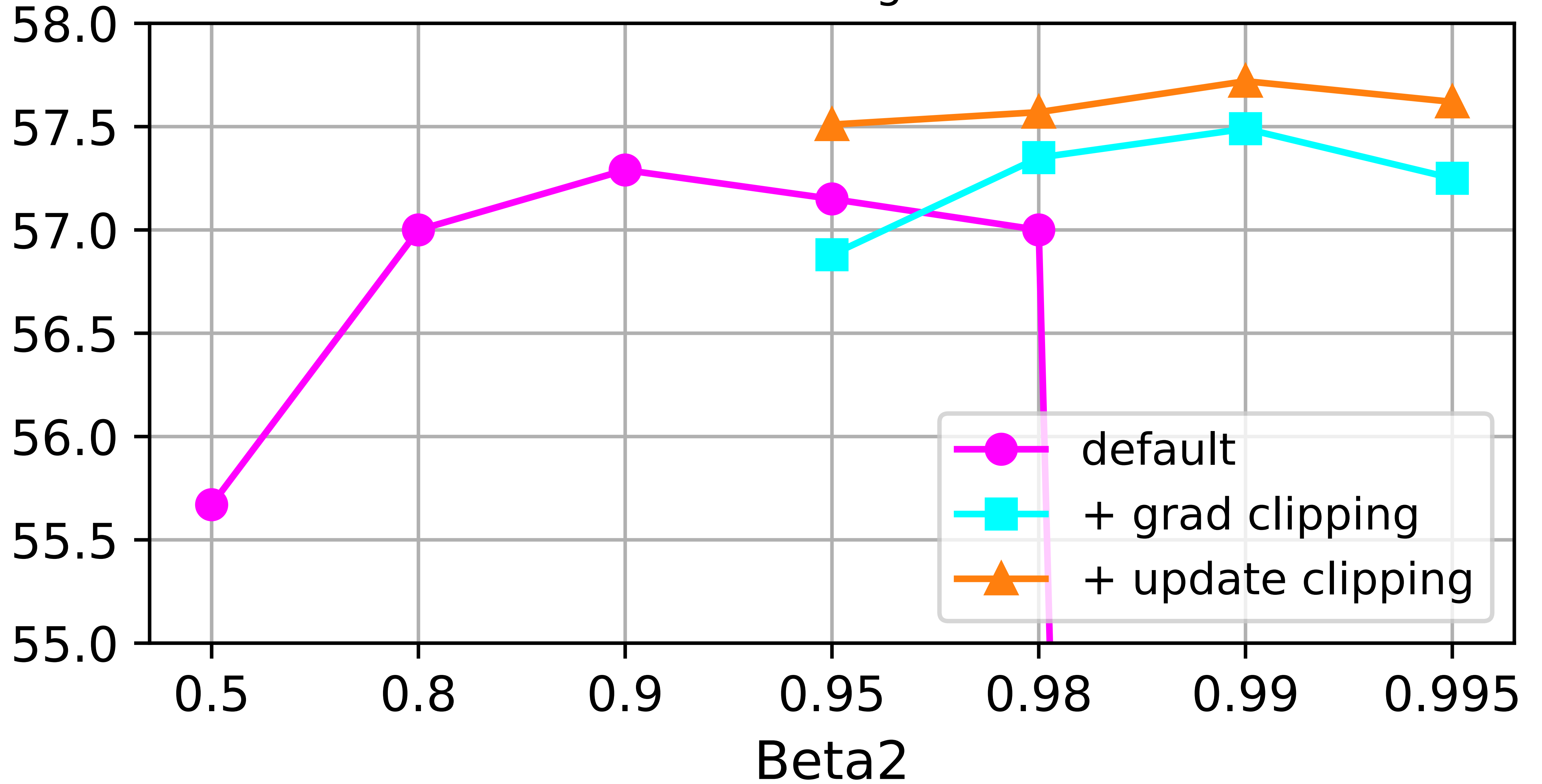
---

ViT-Huge, Beta2 = 0.99



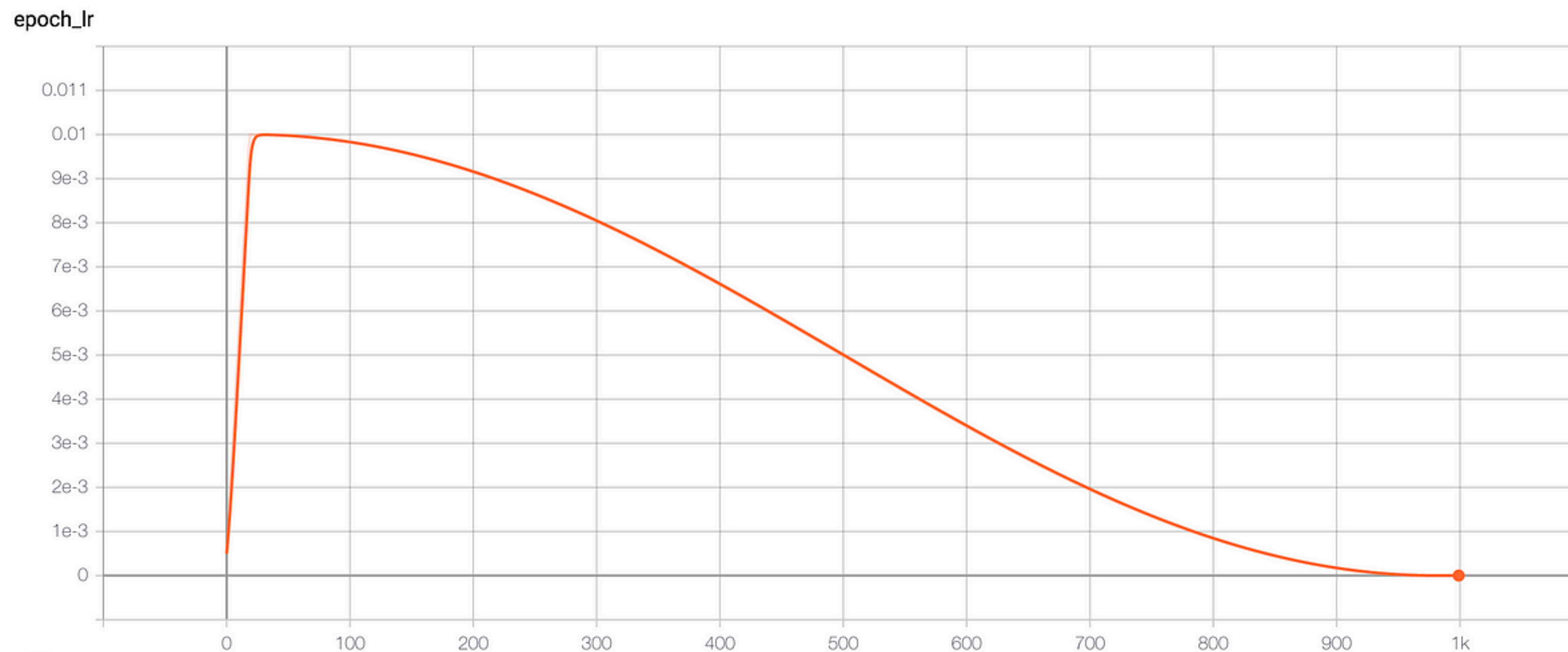
Zero-shot ImageNet accuracy

ViT-Huge



# Other tips & tricks for optimization

- Use around 5000 iterations of linear warm-up for the learning rate
- After that, use cosine-decay



# In things still go wrong...

- Inspect your data
- Use a smaller learning rate
- Using bf16 instead of fp16

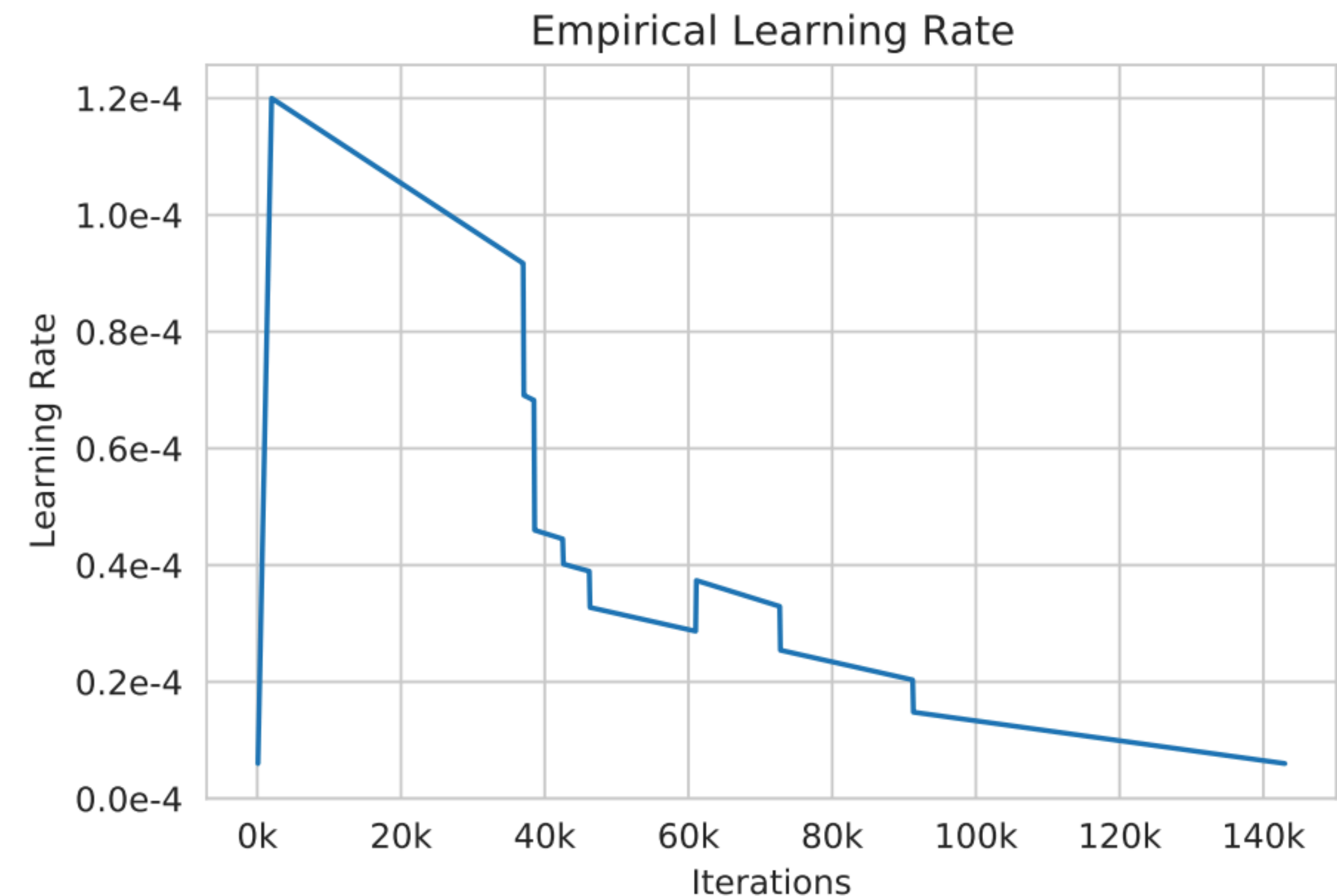


Figure 1: **Empirical LR schedule.** We found that lowering learning rate was helpful for avoiding instabilities.



# What we didn't cover and may be important

- Scaling learning rate and initialization based on width!

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

- There are other fixes for this, including scaling LR by RMS(weights) (see PaLM paper), or “principled” approaches such as mu-transfer

---

# Tensor Programs V: Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer

---

Greg Yang<sup>\*×</sup> Edward J. Hu<sup>\*×†</sup> Igor Babuschkin<sup>°</sup> Szymon Sidor<sup>°</sup> Xiaodong Liu<sup>×</sup>  
David Farhi<sup>°</sup> Nick Ryder<sup>°</sup> Jakub Pachocki<sup>°</sup> Weizhu Chen<sup>×</sup> Jianfeng Gao<sup>×</sup>  
<sup>×</sup>Microsoft Corporation <sup>°</sup>OpenAI

## Abstract

Hyperparameter (HP) tuning in deep learning is an expensive process, prohibitively so for neural networks (NNs) with billions of parameters. We show that, in the recently discovered Maximal Update Parametrization ( $\mu$ P), many optimal HPs remain stable even as model size changes. This leads to a new HP tuning paradigm we call  $\mu$ Transfer: parametrize the target model in  $\mu$ P, tune the HP indirectly on a smaller model, and *zero-shot transfer* them to the full-sized model, i.e., without directly tuning the latter at all. We verify  $\mu$ Transfer on Transformer and ResNet. For example, 1) by transferring pretraining HPs from a model of 13M parameters, we outperform published numbers of BERT-large (350M parameters), with a total tuning cost equivalent to pretraining BERT-large once; 2) by transferring from 40M parameters, we outperform published numbers of the 6.7B GPT-3 model, with tuning cost only 7% of total pretraining cost. A Pytorch implementation of our technique can be found at [github.com/microsoft/mup](https://github.com/microsoft/mup) and installable via pip `install mup`.

