

Meeting — Summary LLM

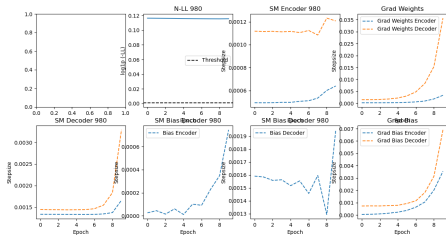
Bas Donkers, Dennis Fok, Finn-Ole Höner

May 13, 2024

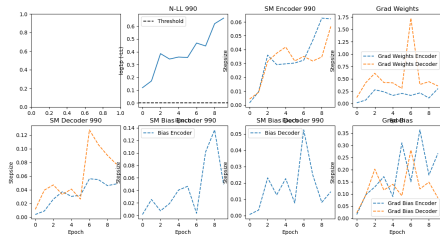
- Gradient spikes first, then large step in decoder / encoder
- These steps seem to affect one hidden unit more strongly than the other
- Increase of Decoder Weight gradients seems to precede the increase of the -LL
- The other gradients seem to react to this jump
- Gradient-Decoder → Step → Other gradients → Step ...
- `torch.nn.utils.clip_grad_norm_(ae.parameters(), max_norm=dMaxGrad)`¹ seems to help. This clips the gradient if the norm exceeds `dMaxGrad`. See here².

¹https://pytorch.org/docs/stable/generated/torch.nn.utils.clip_grad_norm_.html#torch.nn.utils.clip_grad_norm_

²https://www.cs.toronto.edu/~rgrosse/courses/csc321_2017/readings/L15%20Exploding%20and%20Vanishing%20Gradients.pdf#page=6.78

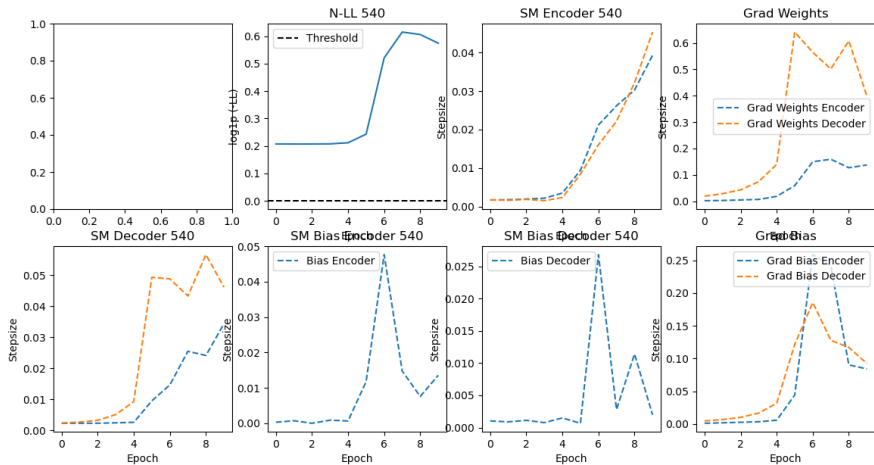


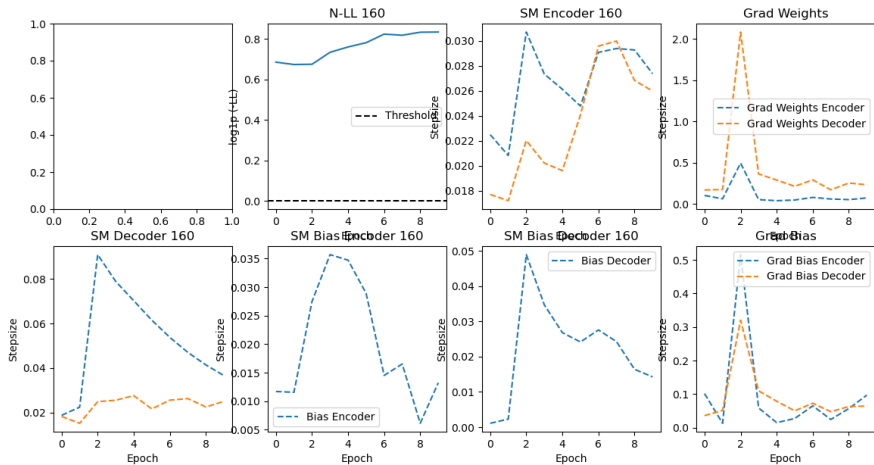
(a) Epoch 970-980

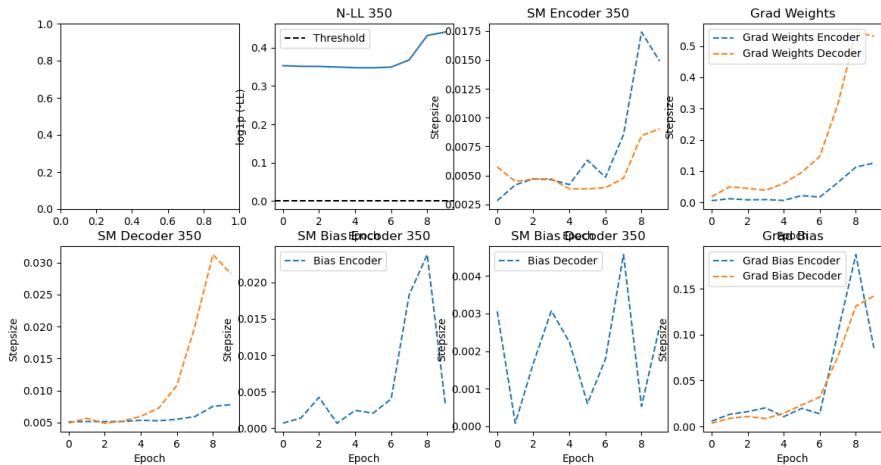


(a) Epoch 980-990

Figure 2: Spikes: Notice how the decoder gradient norm increases first, causing a large step and in-turn affecting the other gradients. The scales of the y-axis differ across all plots.







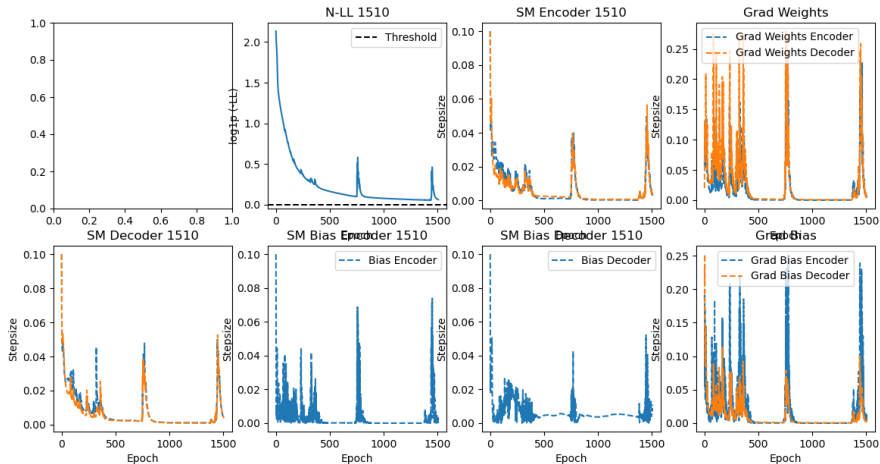


Figure 3: Fixed clipping at 0.3. Later in the optimization this is too large.