

Stochastic Weight Averaging on a Transformer Network

Finding a flatter loss surface for textual data

Elisabeth Stockinger, Askhat Issakov

Introduction: Text Simplification

Text Simplification = **Machine Translation**

Input language: **Original text**

Output language: **Simplified text**

WikiSmall Dataset

Source

August is the eighth month of the year in the Gregorian Calendar and one of seven Gregorian months with the length of 31 days.

Target

August is the eighth month of the year. It has 31 days.

Data Preparation and Loading

Spacy

- English text **tokenization**

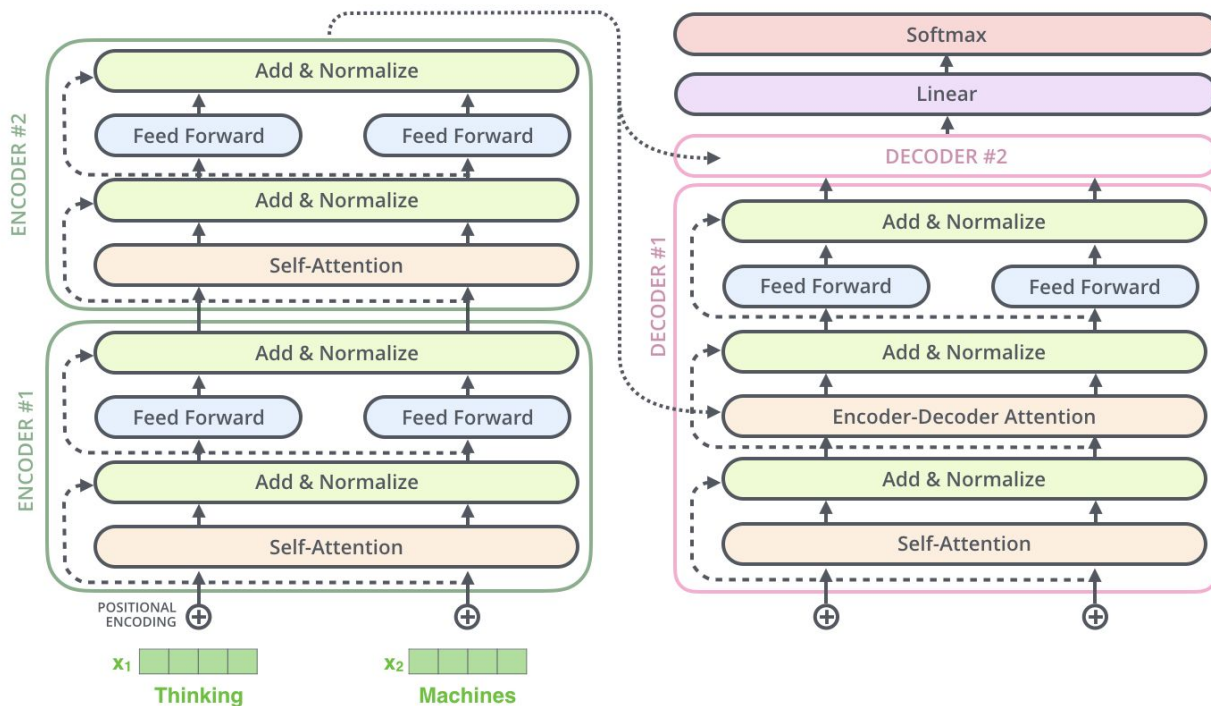
Torchtext

- **Field** - holds the tokenized text
- **TranslationDataset** - associates original with target texts
- **Vocab** - converts text into vectors of numbers

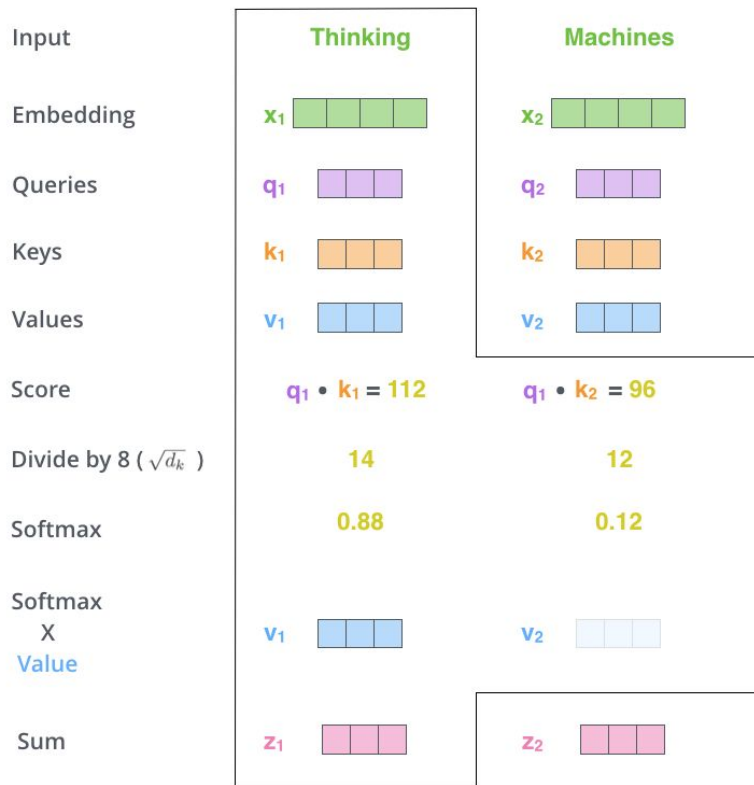
<https://spacy.io/usage/spacy-101>

<https://torchtext.readthedocs.io/en/latest/index.html>

Transformer Network



Self-Attention



Multi-Headed Attention

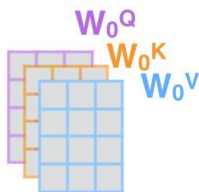
1) This is our input sentence*

Thinking
Machines

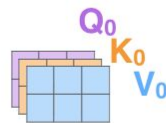
2) We embed each word*



3) Split into 8 heads. We multiply X or R with weight matrices



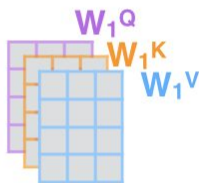
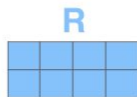
4) Calculate attention using the resulting $Q/K/V$ matrices



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



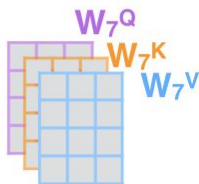
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

...

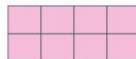
...



W^O

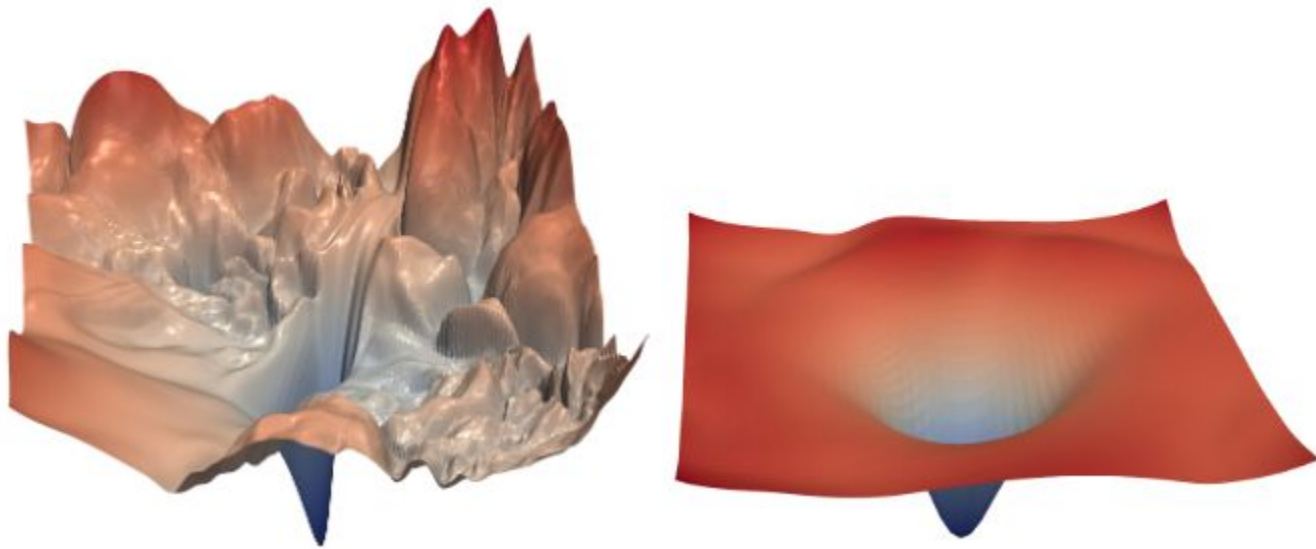


Z



Stochastic Weight Averaging

Understanding Loss Surfaces

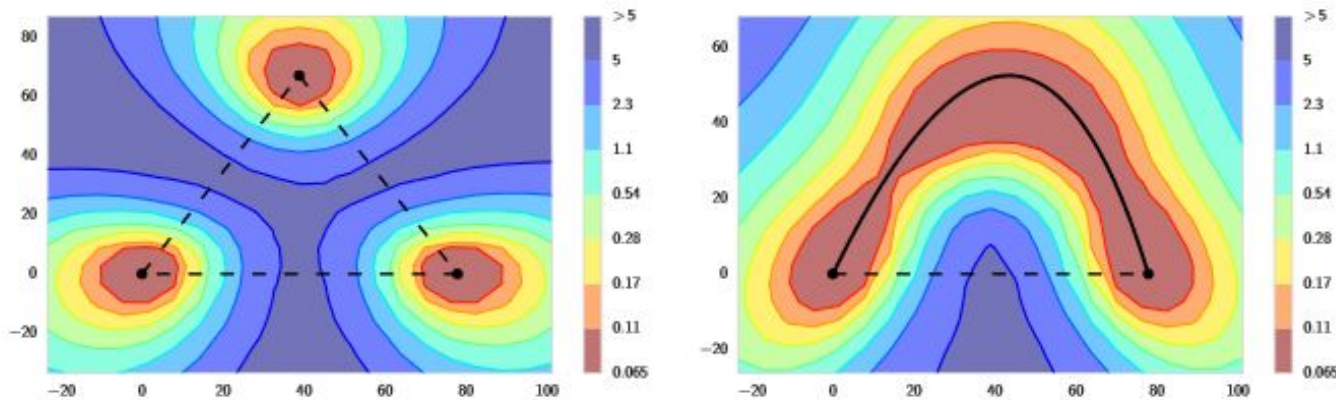


The loss surfaces of ResNet-56 with and without skip connections (visualized in low dimension)

Understanding Loss Surfaces

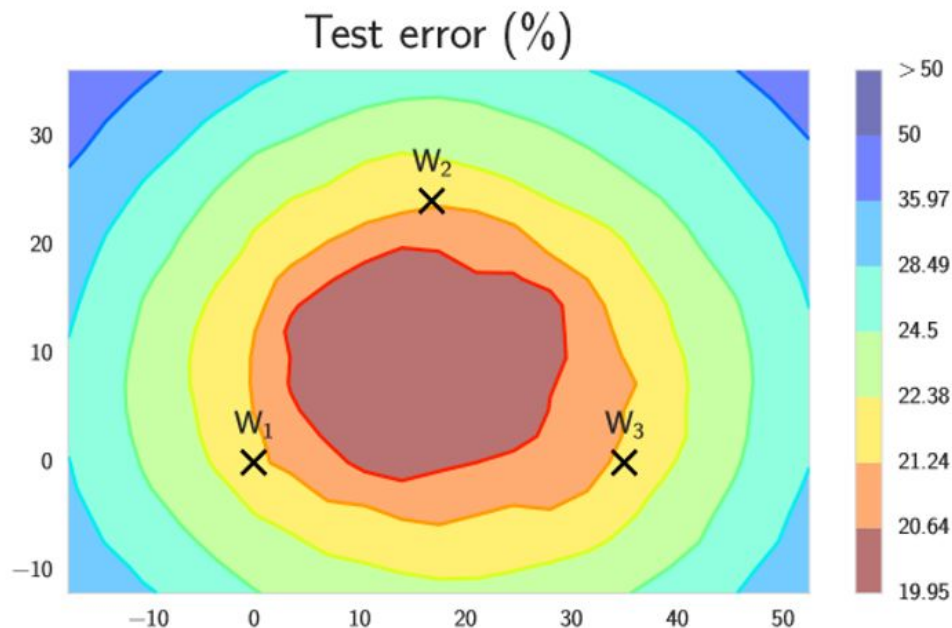
Wide minima generalize better than sharp minima

Local minima can be connected by curves of near-constant loss



The L2-regularized cross-entropy train loss surface of a ResNet-164 on CIFAR-100.

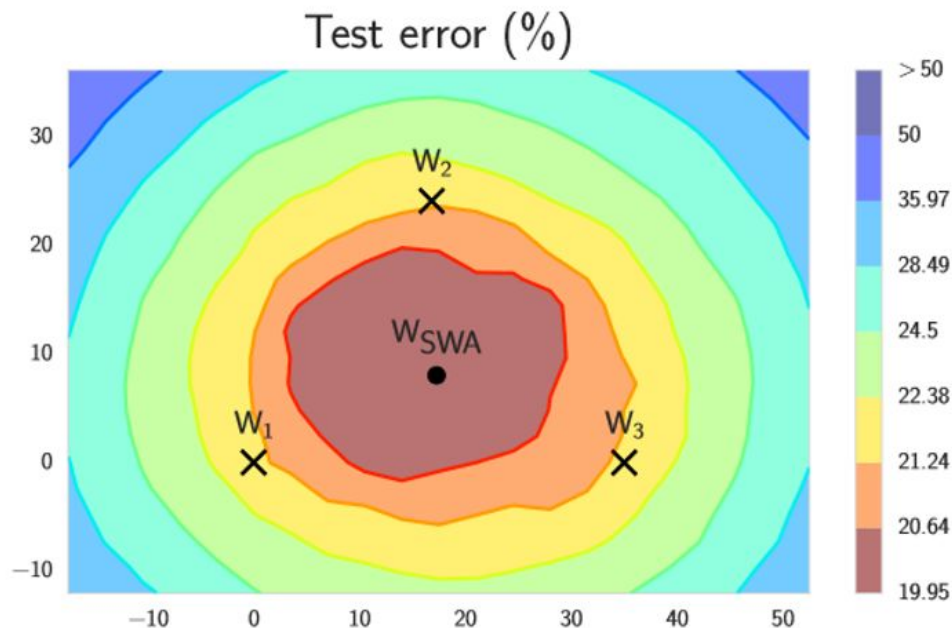
Understanding Loss Surfaces



W_1, W_2 and W_3 are local minima found through SGD

Minima found by SGD are **constrained to the surface of a sphere** of high-dimensional Gaussian.

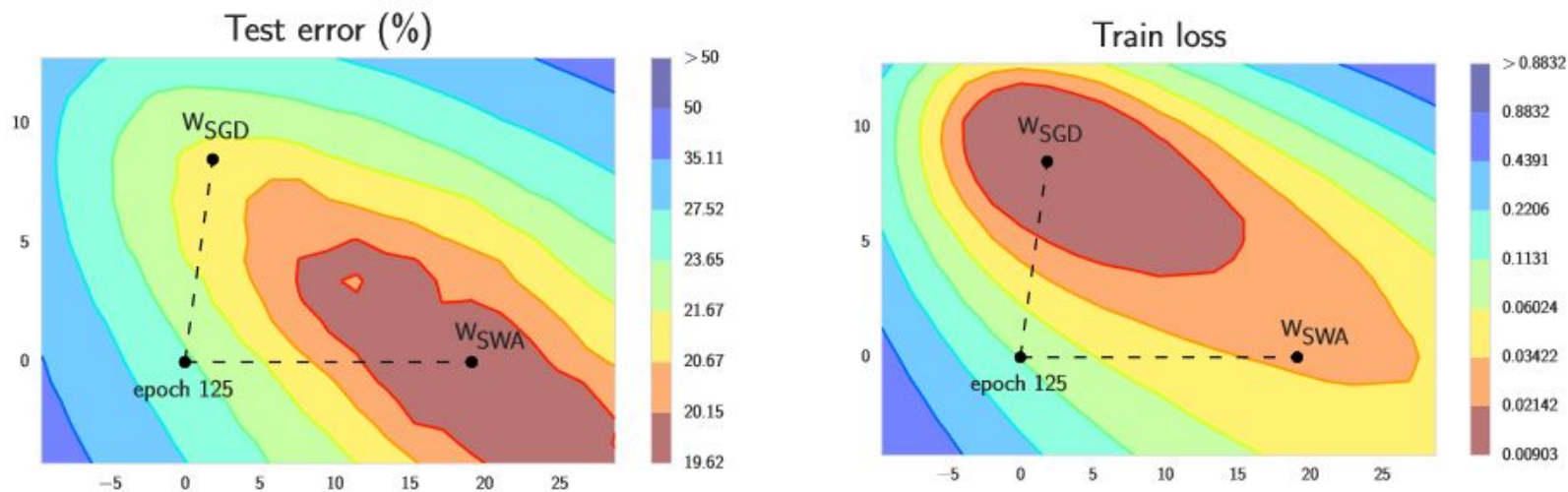
Stochastic Weight Averaging



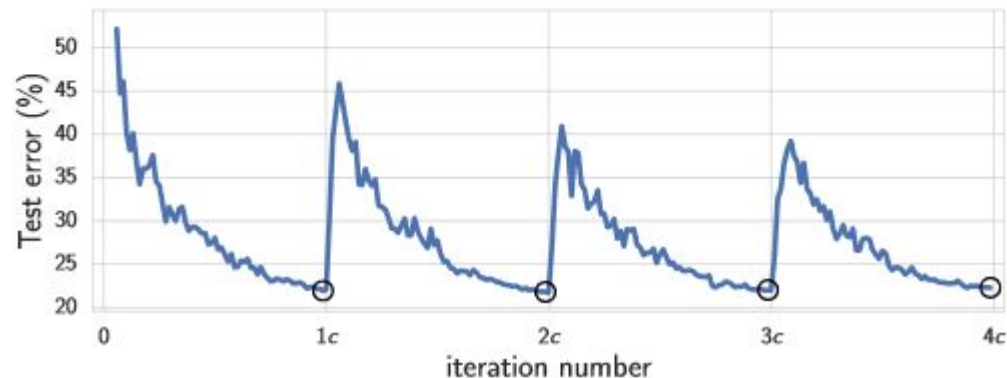
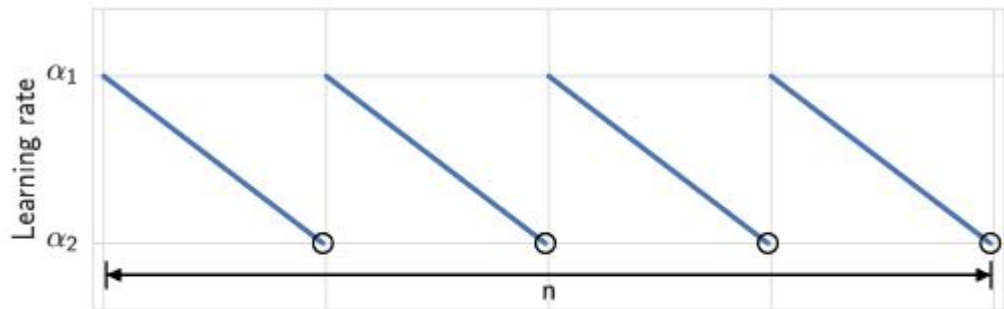
W_1, W_2 and W_3 are local minima found through SGD

SWA lets us **enter the sphere**

Stochastic Weight Averaging



Stochastic Weight Averaging



Cyclical learning rate:
Linearly decrease the learning rate from an upper bound to a lower bound in each cycle.

Stochastic Weight Averaging Algorithm

$w_{SWA}, w \leftarrow \hat{w}$

for $i \leftarrow 1, 2, \dots, n$ **do**

$\alpha \leftarrow \alpha(i)$

$w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$

if $\text{mod}(i, c) = 0$ **then**

$n_{models} \leftarrow i/c$

$w_{SWA} \leftarrow \frac{w_{SWA} \cdot n_{models} + w}{n_{models} + 1}$

end if

end for

Stochastic Weight Averaging Algorithm

$$w_{SWA}, w \leftarrow \hat{w}$$

for $i \leftarrow 1, 2, \dots, n$ **do**

$$\alpha \leftarrow \alpha(i)$$

$$w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$$

if $\text{mod}(i, c) = 0$ **then**

$$n_{\text{models}} \leftarrow i/c$$

$$w_{SWA} \leftarrow \frac{w_{SWA} \cdot n_{\text{models}} + w}{n_{\text{models}} + 1}$$

end if

end for

$$\alpha(i) = (1 - t(i))\alpha_1 + t(i)\alpha_2$$

$$t(i) = \frac{1}{c}(\text{mod}(i - 1, c) + 1)$$

Results

Train Network for 60 Epochs using SGD

SGD

SWA, CL 5, LR 0.001 – 0.1

SWA, CL 5, LR 0.05 – 0.1

SWA, CL 10, LR 0.001 – 0.1

SWA, CL 10, LR 0.05 – 0.1

SWA, CL 10, LR 0.05 – 0.1

SWA, constant LR 0.01

Results

Train Network for 60 Epochs using SGD		SGD	1.4296
1.6330		SWA, CL 5, LR 0.001 – 0.1	1.0845
		SWA, CL 5, LR 0.05 – 0.1	1.1580
		SWA, CL 10, LR 0.001 – 0.1	1.1522
		SWA, CL 10, LR 0.05 – 0.1	1.2842
		SWA, constant LR 0.01	1.3986

Results

SGD after 1 budget	1.6330
---------------------------	---------------

SGD after 1 budget and 15 epochs	1.4296
---	---------------

After 15 epochs of SWA:

Cycle Length	Learning Rate	
5	0.001 - 0.01	1.0845
	0.05 - 0.01	1.1580
10	0.001 - 0.01	1.1522
	0.05 - 0.01	1.2842
Constant	0.01	1.3986

Translated samples

Source

while in england hendrix invited cox to join him in a new band cox declines preferring to work in various backing bands .

Reference

while in england hendrix asked cox to join him in a new band cox said no .

Output

it is a of is a in of the and is a of .

Translated samples

Source

it was founded in the 14th century by genoese colonists , who employed large numbers of workmen (calfats) in repairing ships .

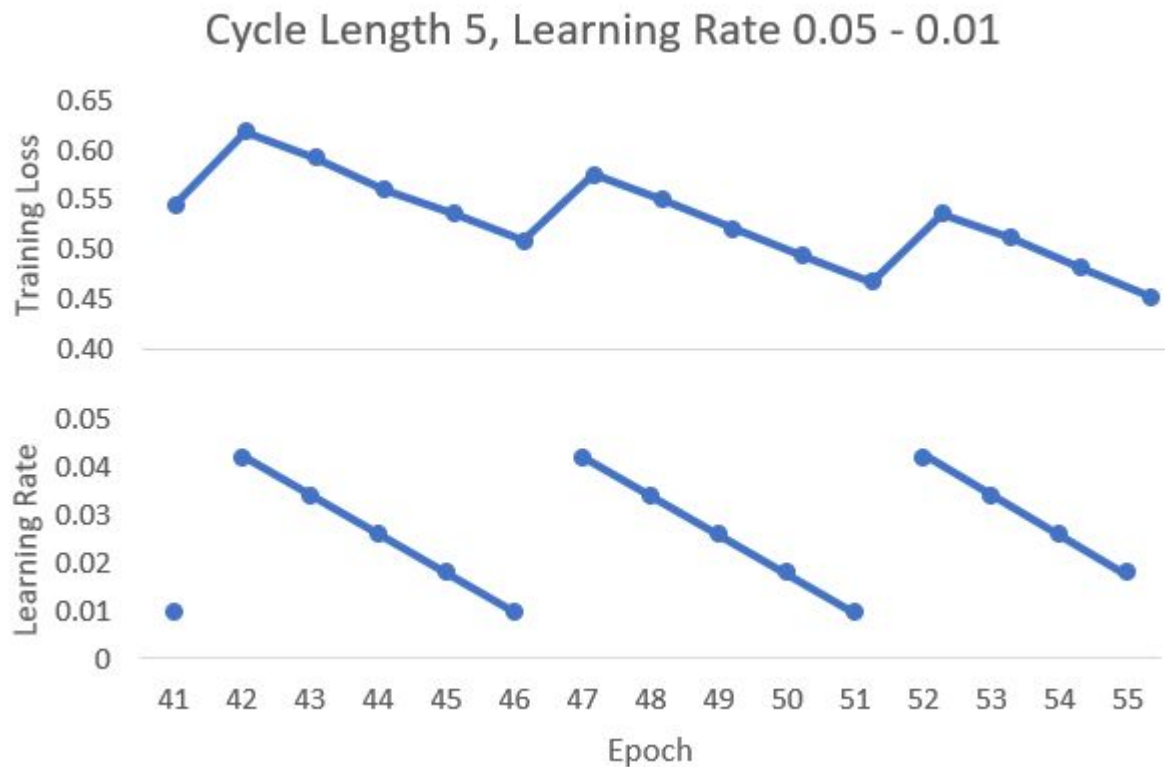
Reference

calafat was started in the 14th century by genoese colonists .

Output

chinese successor following person muslim revolution god son saint house
historian meeting dynasty roman patron professor rare prayer lance

Loss VS epoch and learning rate



Challenges

- Dataset size - Colab instance could not load it into memory
- GPU - Colab instance runs out of GPU memory while training
- Training time - very large training time required