

Probability of the entire token sequence is therefore given by Equation 4.

$$\mathbf{p}_t : \{1, \dots, V\} \rightarrow [0, 1] \quad ; \quad \mathbf{Y}_t \sim \mathbf{p}_t \quad (2)$$

$$\mathbf{p}_t(\mathbf{y}_t) := \mathbb{P}(\mathbf{Y}_t = \mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{I}) \quad (3)$$

$$\mathbb{P}(\mathbf{y} | \mathbf{I}) = \prod_{t=1}^{\tau} \mathbf{p}_t(\mathbf{y}_t) \quad (4)$$

As is typical with sequence generators, the training objective here is to maximize probability of the target sequence \mathbf{y}^{GT} . We use the standard per-word cross-entropy objective (Equation 5), modified slightly for the mini-batch (Equation 6). We did not use any regularization objective, relying instead on dropout, data-augmentations and synthetic data to provide regularization.

$$\mathcal{L}_{seq} = -\frac{1}{\tau} \sum_t \ln(\mathbf{p}_t(\mathbf{y}_t^{GT})) \quad ; \quad \tau \equiv \text{sequence length} \quad (5)$$

$$\mathcal{L}_{batch} = -\frac{1}{n} \sum_{batch} \sum_t \ln(\mathbf{p}_t(\mathbf{y}_t^{GT})) \quad ; \quad n \equiv \# \text{ of tokens in batch} \quad (6)$$

The final Linear layer of the decoder (Figure. 3a) is a 1x1 convolution function that produces logits which are then normalized by softmax to produce \mathbf{p}_t .

Combination of Vision and NLP One of the strengths of our architecture is in the combination of Vision and Language models. CNNs such as ResNet are considered best for processing image data. And Transformers are considered best for Language Modeling (LM) and Natural Language Understanding (NLU) tasks [7, 24, 25], possessing properties that are very useful in dealing with noisy and incomplete text that often occurs in real handwriting. Having both the visual feature map and a language model, the model can do a much better job than one relying on visual features alone.

Inference We use simple greedy decoding, which picks the highest probability token at each step. Beam search decoding [8] did not yield any accuracy improvement indicating that the model is quite opinionated / confident.

4 Training Configuration and Procedure

The base configuration uses grayscale images scaled down to 140-150 dots per inch. Higher resolutions yielded slightly better accuracy at the cost of compute and memory. We use the 34-layer configuration of ResNet, but have also successfully trained the 18-layer and 50-layer configurations; larger models tending to do better in general as expected.

The following is the base configuration of the Transformer stack:

- N (number of layers) = 6