

Supplementary Material

This document provides supplementary material for the paper “Text and Style Conditioned GAN for Generation of Offline-Handwriting Lines” submitted to BMVC 2020, including details of the human study described in the paper, additional image results, additional experimental ablation study results, and architectural details for the networks described in the paper. The sections are as follows:

- **S.1** Details on FID (and GS) computation.
- **S.2** Details on human experiment.
- **S.3** Additional generation results.
- **S.4** Additional ablation results.
- **S.5** Network specifications of each model part.

S.1 Discussion of FID evaluation and GS details

FID [12] is evaluated by passing an image through the convolutional network Inception-v3 and computing statistics on the average pooled features. Inception-v3 was designed to accept images of size 299×299 , and thus most implementations of FID rescale images to this size before feeding them to the network. In most situations this is fine since GANs typically generate square images. However, in the case of handwriting, particularly lines, images are generally much wider than they are tall. Resizing them to be square causes significant distortions to the image. Thus, it would make sense to resize images to a height of 299 and maintain the aspect ratio. Since Inception-v3 is fully convolutional up to the average pooling, it can accept variable sized images. We evaluated FID with both the original square resizing and aspect ratio preserving resizing. We found the scores produced when preserving the aspect ratio appeared closest to the FID reported in [2] and [6] and thus assume these authors applied something similar, although they do not report this. We follow [2] in using 25,000 training set images and generate 25,000 images using the same lexicon (words or lines depending on dataset), but styles extracted from the test set. Like [6], we only run the experiment once.

When comparing our generated images to RIMES words, there is a distribution difference caused by segmentation differences. RIMES words are segmented tightly to each word. Our model is trained on RIMES lines, which generally have more whitespace on the top and bottom of each word. Fig. 5 demonstrates this difference. To make comparison more fair, we crop our generated words on the top and bottom to the first ink pixel (value less than 200). This cropping resembles the segmentation of the word images and slightly improves our FID score.

We also question in general the validity of using FID score for handwriting images. As Inception-v3 is trained on natural images, not handwriting, FID seems ill-suited for evaluating the quality of handwriting images. Further investigation is required into the topic of applying FID to image domains other than natural images.

For GS [22], the data is expected to all be the same size. Because the dataset has variable width images and our method produces variable width images, we pad images to be the same width. Neither [2] nor [6] report how they handled this. Like [6], we only run the experiment once.

S.2 Human Study Details

We submitted 78 image tasks to Amazon Mechanical Turk (35 real, 35 generated, 8 poorly generated), requesting 200 workers to review each image. Each task consisted of instructions, with example images, a task image (real, generated, or poorly generated) and two multiple choice questions. The first question asked the worker to select the correct transcription for the task image. Two choices were shown, one with the correct transcription, the other a permutation of the correct transcription’s words (where the first and last words remained in the same place). We removed punctuation so the permutation didn’t create artifacts that made the choice too easy. This was to ensure the worker actually looked at the image and was paying attention to what they were doing. The second asked if they thought the image was written by a human or a computer.

The interface the workers saw can be seen in Fig. S1. The correct and incorrect transcription options were randomly ordered, the options between human and computer remained in the same order.

The real instances used in the study were randomly selected from the test set. The generated images used the same text as the selected real instances, but the styles were from interpolations between styles extracted from randomly selected test set images.

To help measure the reliability of the workers, we included poorly generated images which should appear to not be written by a human. These were created using a model only trained 2,000 iterations. The responses on these images were not included in the final evaluation, but were held out to help gauge the confidence that can be placed in the workers efforts. The poorly generated images used in the study are shown in Fig. S2. The generated and dataset images used in the study are in Figs. S3 and S4 respectively.

The transcription question was used to filter out workers which were unreliable (likely clicking random responses to complete the tasks quickly). We only used workers who had at least 90% accuracy on transcription (permutations can sometimes be very close to the correct transcription leading to some error in even engaged workers). Additionally, we only used workers we had at least 6 responses for. The selected workers had 89.5% accuracy on the poorly generated images, the left-out workers had 79.0% accuracy.

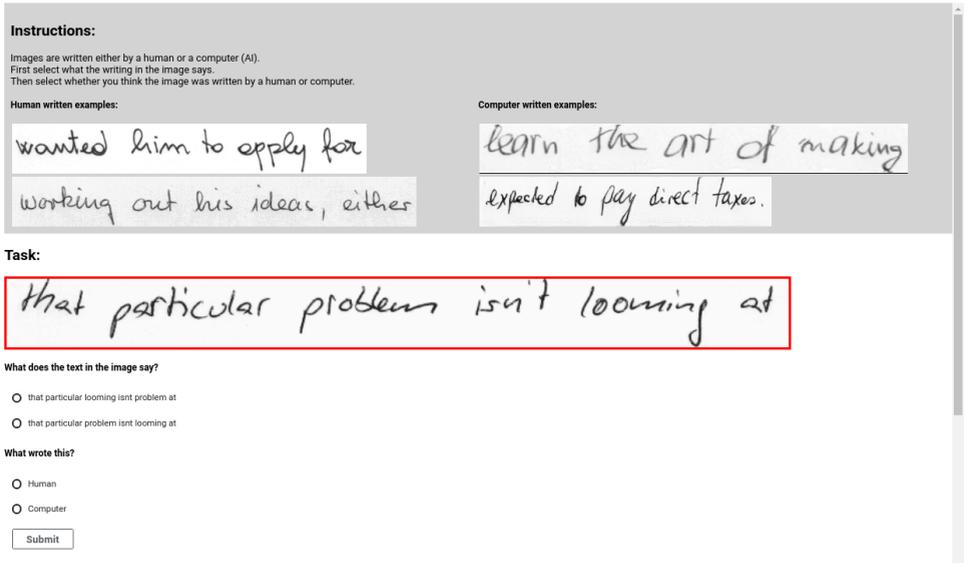


Figure S1: A screenshot of the interface the workers saw when completing a task. The example images remained the same each task. The order in which the correct and incorrect transcription responses were placed was random. We kept the task image large so detail could be seen.

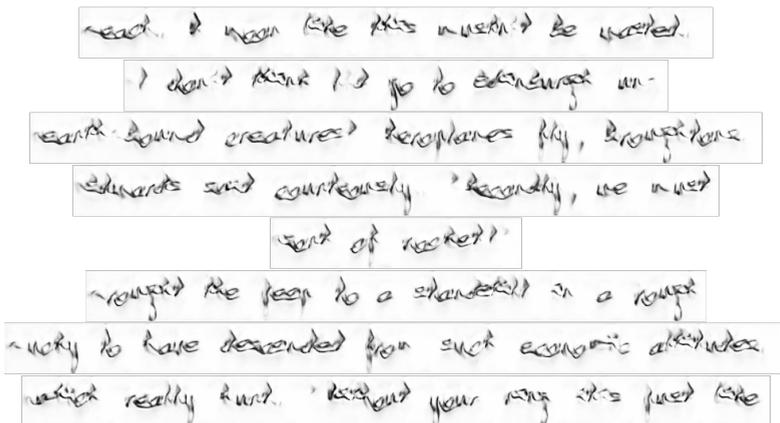


Figure S2: Poorly generated images from an intentionally under-trained model used in human study to evaluate participant ability or attention. These samples are *not* from our final model.

after Simone had left expect her to accept his
 we were to go no further unless and
 opportunity presented itself. He must
 I'm such a dull fellow, really." "Dull?" She
 With an air of resignation he sent Judy, his
 his toes. He stretched his arms
 to other day, my lord. A being of whom
 disposed towards each other.
 particular surgery?" "You are a fool,
 and have a look at the mining camp.
 back there on a great plan we have.' 'Go back?'
 Mr. Copthorne was on dry land in a church
 got all these lovely things" - she waved a
 Something had cropped up which required Nigel's attention,
 A single-decker," he elaborated. Daggers, the
 R?Bs as 'Seaweed', and a youngish, sharp-eyed
 him at the temporary bridge over the
 the housemen think of me as a
 allies being more inferior than formerly to W.C.U. in the
 that particular problem isn't looming at
 ringing of a doorbell was to him
 sitting just inside having coffee.
 radar glowed. Occasionally minute spots
 hell lying there, knowing she was
 Gavin and the girl who had got
 I'm such a dull fellow, really." "Dull?" She
 her word. There was enough evidence,
 I don't think I'd go to Edinburgh un-
 who is tough enough to change him."
 had tin mugs filled with hot black coffee
 a substantial breakfast. Although usually a very
 earth-bound creatures? Aeroplanes fly, Broughtons.
 over a fairly wide area working singly or
 found us a couple of boxes to sit on.
 Haris had trampled on his. It was unthinkable!

Figure S3: Generated images used in human study that were generated using random styles (i.e. random interpolation of style vectors extracted from random pairs of real images from IAM) and random text from the IAM corpus.

allis being more inferior than ~~found~~ to W. C. U. in the
 "I'm such a dull fellow, really." "Dull?" she
 over a fairly wide area working only on
 A single - decker," he elaborated. Daggars, the
 found us a couple of boxes to sit on.
 With an air of resignation he sent Judy, his
 her word. There was enough evidence,
 got all these lovely things'-she naved a
 to ~~other~~ other day, my land. A being of whom
 Mr. Cophorne was on dry land in a duds ~~at~~
 and have a look at the mining camp.
 ringing of a doorbell was to him
 his toes. He stretched his arms
 I don't think I'd go to Edinburgh in-
 opportunity presented itself. He must
 "I'm such a dull fellow, really." "Dull?" she
 hell lying there, knowing she was
 who is tough enough to change him."
 *?2s as 'Seaweed', and a youngish, sharp-eyed
 disposed towards one another.
 that particular problem isn't looming at
 had tin mugs filled with hot black coffee
 Gavin and the girl who had got
 the housemen think of me as a
 something had cropped up which required Nigel's attention,
 him at the temporary bridge over the
 sitting just inside having coffee.
 earth-bound creatures? Heroplaus fly, Broughtons
 we were to go no further unless and
 particular surgesy?" "You are a fool,
 back there on a great plan we have.' 'Go back?'
 a substantial breakfast. Although usually a very
 radar glowed. occasionally minute spots
 after Simone had left expect her to accept his
 Maris had trampled on his. It was unthinkable!

Figure S4: Dataset images used in human study. These are randomly sampled from IAM.

S.3 Additional Generation Results

We here show additional results from our model. Fig. S5 shows additional examples of style interpolation. Figs. S6 and S7 shows generation using random interpolated/extrapolated styles with fixed and varying text respectively. Figs. S8 and S9 show reconstruction results.

panions he greeted courteously by name as they entered
 will shorten the stirrups up to the saddle-skirts, and
 g the ruined walls and paving-stones of an ancient high
 nd. I will answer some of your questions, if that will
 any - yours not least. In any case we did not kill h
 r despair is only for those who see the end beyond all.
 blades were forged many long years ago by Men of West
 them. Frodo has a better head for that sort of thin
 n-king. Of him the harpers sadly sing: the last whose s
 rd, and harps of gold they brought to him. They cloth
 ever you return? 'Not this at least. said Boromir
 in his chair, and looked at the farmer with an unfrien
 he Lady Galadriel approaching. Tall and white and fair s
 d wildly. They were too overjoyed to hear him speak t
 the Shire. The Sackville-Bagginses were not forgotte
 ths, every Baggins, Boffin, Took, Brandybuck, Grubb, C
 country. There the River flows in stony vale amid hig
 aiting perhaps for a change of days, and he will not st
 But either Gandalf was astray, or else the land had
 any - tiered branches and amid their ever-moving leave
 terror they bore their riders into the rushing flood
 them. Bilbo had not much to say of himself. When he h
 get away without those cursed goblins seeing us.' 'Pe
 , for it is a healing plant that the Men of the West br
 understand, say, a Dwarf, or an Ore, or even an Elf
 he landlord, pausing and snapping his fingers. 'Ah, yes
 quite and ambiguous. It is also false, though naturally
 st to please me, I think; for, of course, they aren't
 ock, pierced by a dark arch like a great gate. It se
 'thril! I have never seen or heard tell of one so fair.
 nd his bootless foot is lasting lame; But Troll don't s
 his wish and need, but especially that one of the lit
 looked over the land ahead, and called to Pippin. '
 g for breath. He saw as through a mist a wide flat cir
 shore. The sunlight glittering on the water dazzled hi
 into tears. Chapter 3. Three is Company 'You ought to
 raight as he could over the wild lands to Weather-top H
 swooned he caught, as through a swirling mist, a glimp

Figure S7: Additional generation results using random extra/interpolations between test set styles using varying text.

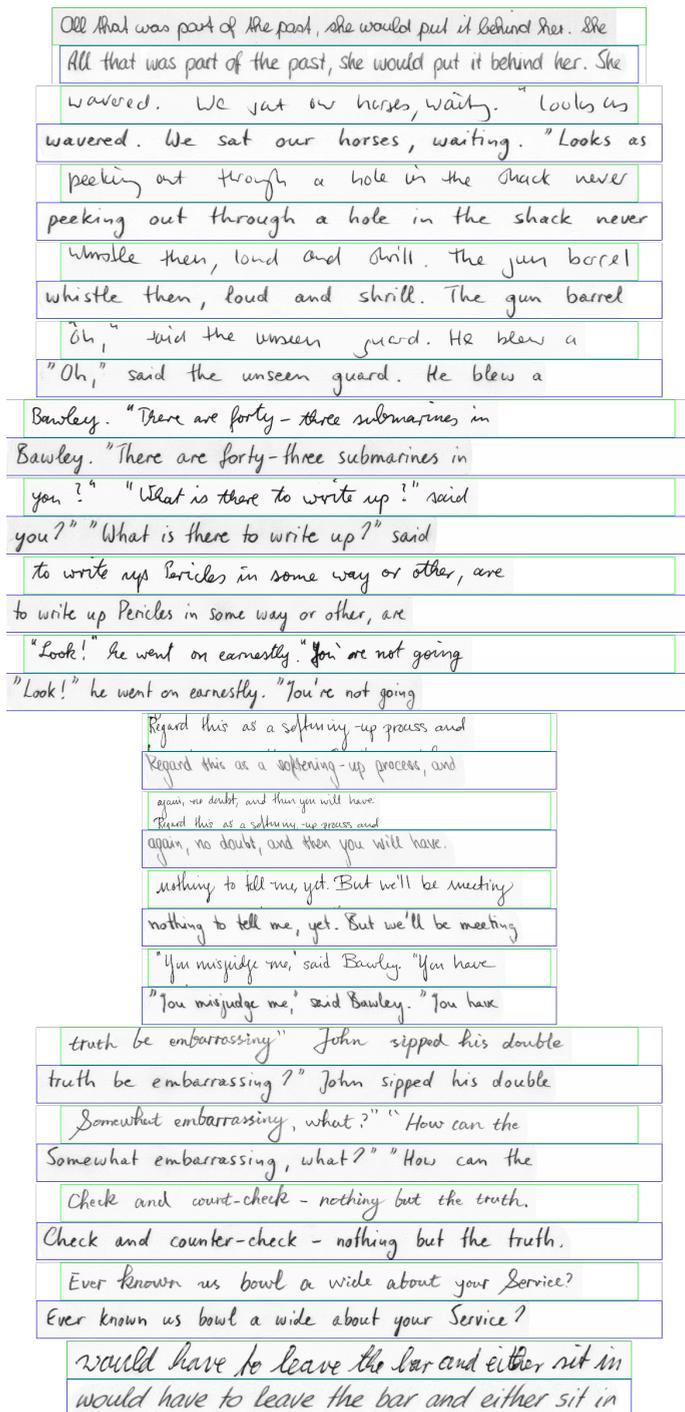


Figure S8: Additional Reconstruction results. Green is original, blue is our model's reconstruction.

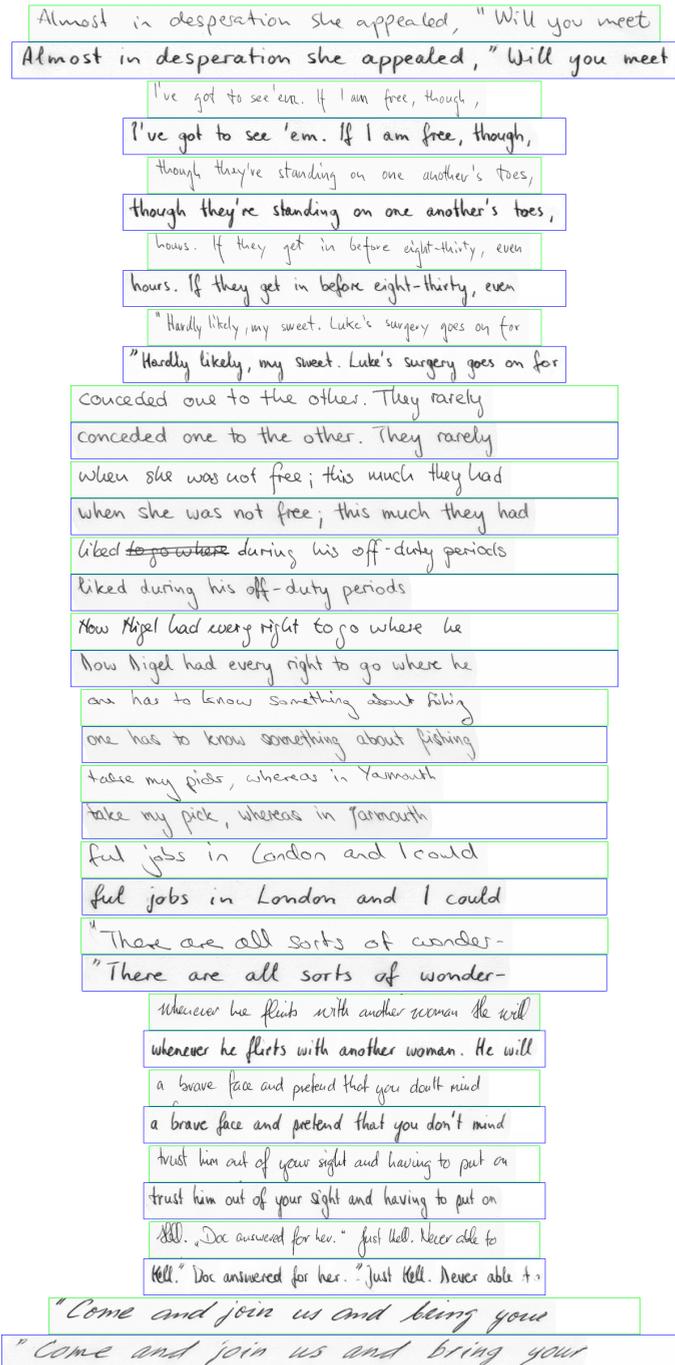


Figure S9: Additional Reconstruction results. Green is original, blue is our model's reconstruction.

S.4 Additional Ablation Results

We present additional results for each of the ablation models:

- Fig. S10: No reconstruction loss
- Fig. S11: No adversarial loss
- Fig. S12: No handwriting recognition supervision
- Fig. S13: No character specific components of S
- Fig. S14: No pixel reconstruction loss



Figure S10: Additional ablation results, without the reconstruction losses (random styles).

"Yes," I said. "Something happened all right."

"Yes," I said. "Something happened all right."

Mr. Louis Robbs, landlord of the Traveller's Inn, was sitting up the

Mr. Louis Robbs, landlord of the Traveller's Inn, was sitting up the

more. The first intimation that all was not well came when a

came. The first intimation that all was not well came when a

on these fifteen miles of mountain roads, it disappeared from the earth.

on these fifteen miles of mountain roads, it disappeared from the earth.

So the bus set out for Glasgow. But it never reached that Somerset,

So the bus set out for Glasgow. But it never reached that Somerset,

to entertain for an instant the idea, the

to entertain for an instant the idea, the

myself? But was it so? I allowed myself

myself? But was it so? I allowed myself

Mr Septimus stood in the same case as

Mr Septimus stood in the same case as

Sally and of course Mrs Septimus, for surely

Sally and of course Mrs Septimus, for surely

gazed at recently; and I have no

gazed at recently; and I have no

resemblance to the shrunk heads we had

resemblance to the shrunk heads we had

with black, matted hair and striking

with black, matted hair and striking

Small, repulsive creatures they were,

Small, repulsive creatures they were,

of free men everywhere, could only repel them.

of free men everywhere, could only repel them.

quise which, far from attracting the allegiance

quise which, far from attracting the allegiance

like a masquerade of business interests in dis-

like a masquerade of business interests in dis-

The political life of Aristotle looked more and more

The political life of Aristotle looked more and more

nor sound of flying saucers. So they

nor sound of flying saucers. So they

Figure S12: Additional ablation results, without handwriting recognition supervision.

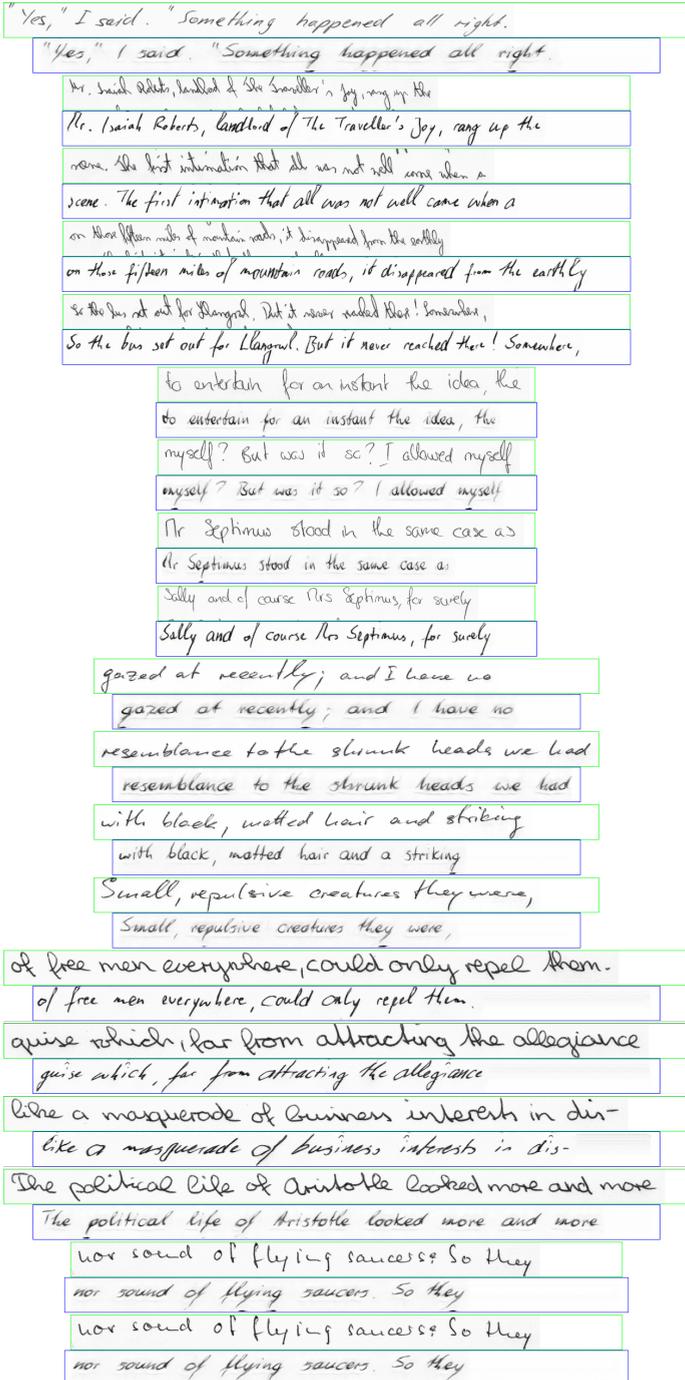


Figure S13: Additional ablation results, without character specific components of S.

"Yes," I said. "Something happened all right."

"Yes," I said. "Something happened all right."

Mr. Smith Roberts, husband of Mrs. Smaller's, jay, rang up the

Dr. Roman Roberts, husband of The Traveler's Joy, rang up the

ness. The first intimation that all was not well came when a

on those fifteen miles of mountain roads, it disappeared from the earth

So the bus set out for Langrath. But it never reached there! Somewhere,

to entertain for an instant the idea, the

to entertain for an instant the idea, the

myself? But was it so? I allowed myself

myself? But was it so? I allowed myself

Mr. Septimus stood in the same case as

Mr. Septimus stood in the same case as

Sally and of course Mrs. Septimus, for surely

Sally and of course Mrs. Septimus, for surely

gazed at recently; and I have no

gazed at recently; and I have no

resemblance to the shrunken heads we had

resemblance to the shrunken heads we had

with black, matted hair and striking

with black, matted hair and a striking

Small, repulsive creatures they were,

Small, repulsive creatures they were,

of free men everywhere, could only repel them.

of free men everywhere, could only repel them.

quise which, far from attracting the allegiance

quise which, far from attracting the allegiance

like a masquerade of business interests in dis-

like a masquerade of business interests in dis-

The political life of Aristotle looked more and more

The political life of Aristotle looked more and more

nor sound of flying saucers? So they

nor sound of flying saucers. So they

Figure S14: Additional ablation results, without pixel-wise reconstruction loss.

S.5 Model Specifications

We present here detailed diagrams of various parts of the model:

- Fig. S15: The handwriting recognition model R
- Fig. S16: The generator G
- Fig. S17: The auxiliary spacing network C
- Fig. S18: The discriminator D
- Fig. S19: The encoder E
- Fig. S20: The style extractor S

The encoder E is trained using the same IAM training set. It is jointly trained with a decoder as an autoencoder with an L1 reconstruction loss and as a handwriting recognition network with a recognition head using the CTC loss. It is trained with the Adam optimizer 6000 iterations with a learning rate of 0.0002.

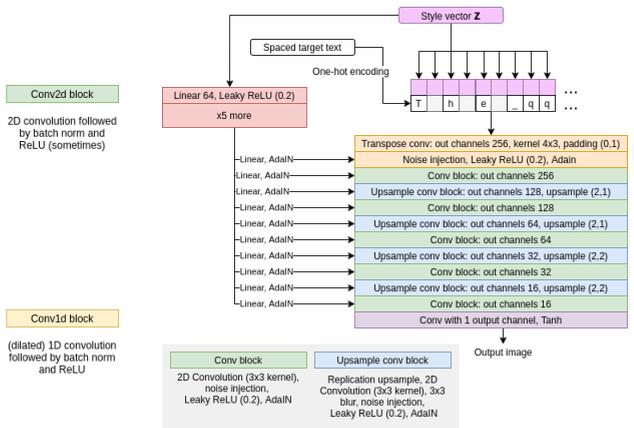
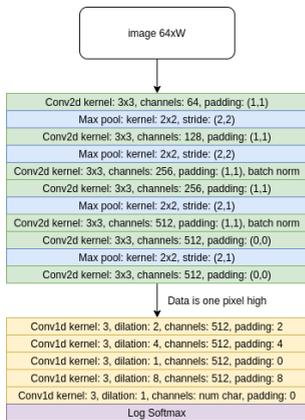


Figure S15: Handwriting recognition network R architecture

Figure S16: Generator G architecture

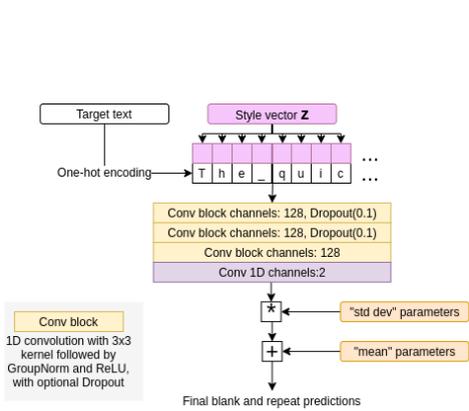


Figure S17: Spacer network C which predicts the spaced text. It predicts the number of blanks preceding each character and the number of times the character should be repeated.

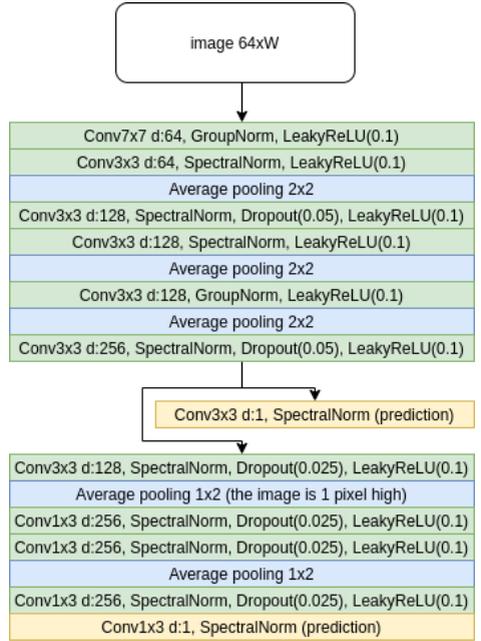


Figure S18: Discriminator D architecture.

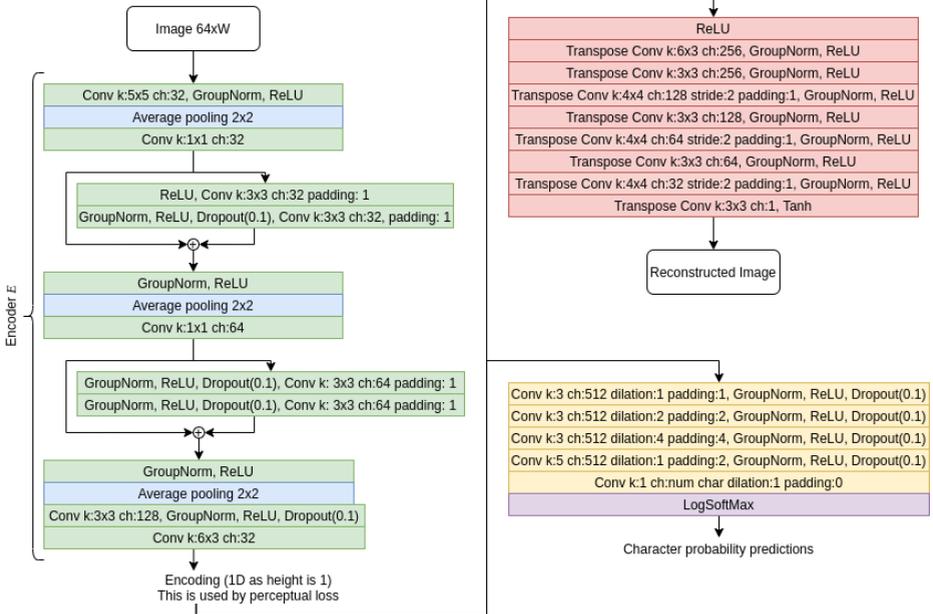


Figure S19: Encoder network E (green) and auxiliary decoder (red) and recognition head (yellow) used to train E .

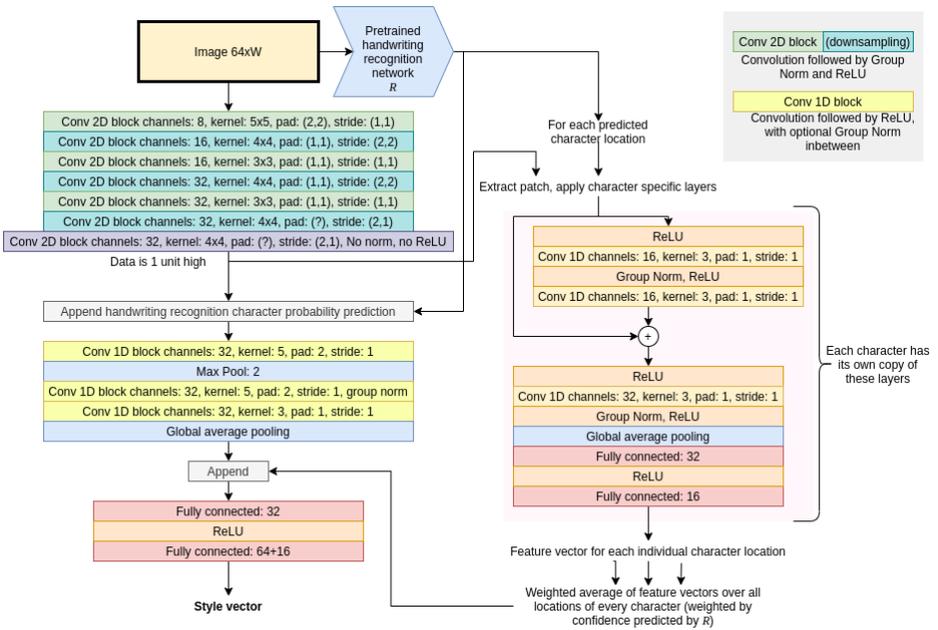


Figure S20: Style Extractor S . It leverages the output of R both as additional input and to (roughly) locate characters. The locations are used to crop features to pass to character specific layers (the learn to extract features for one character).