Images are self-made, part of the IAM Dataset, or credited in the slides.

# Horizontally Variable OCR

# Horizontally Variable OCR
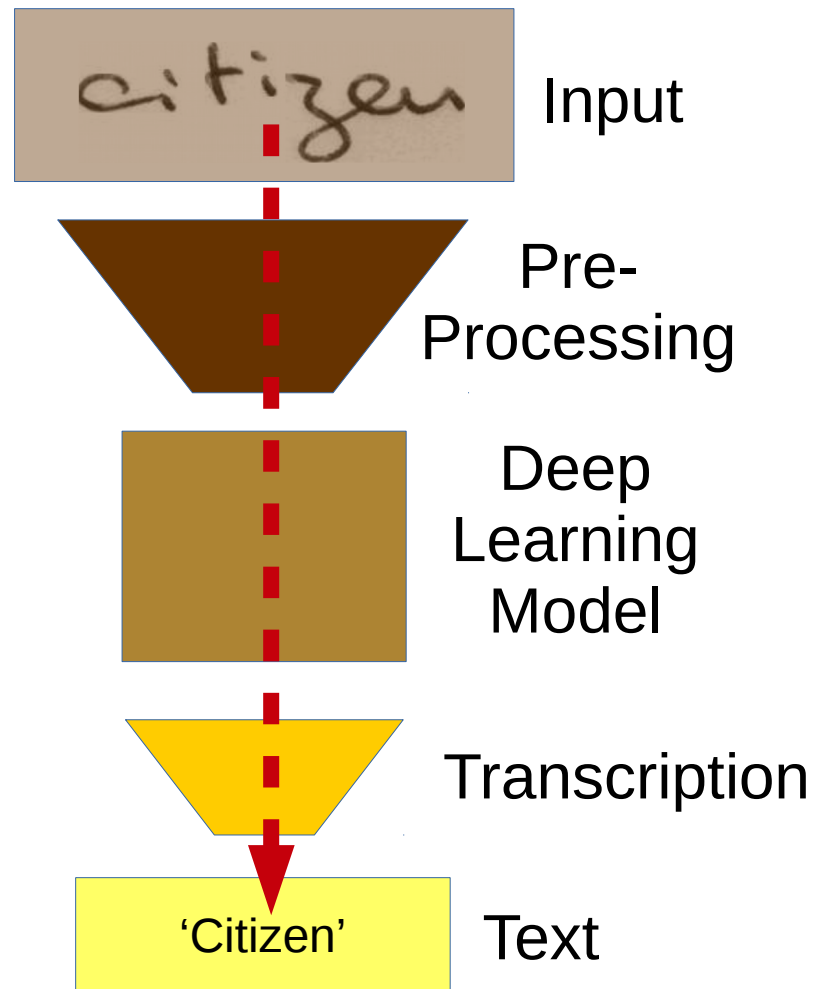
Or something….

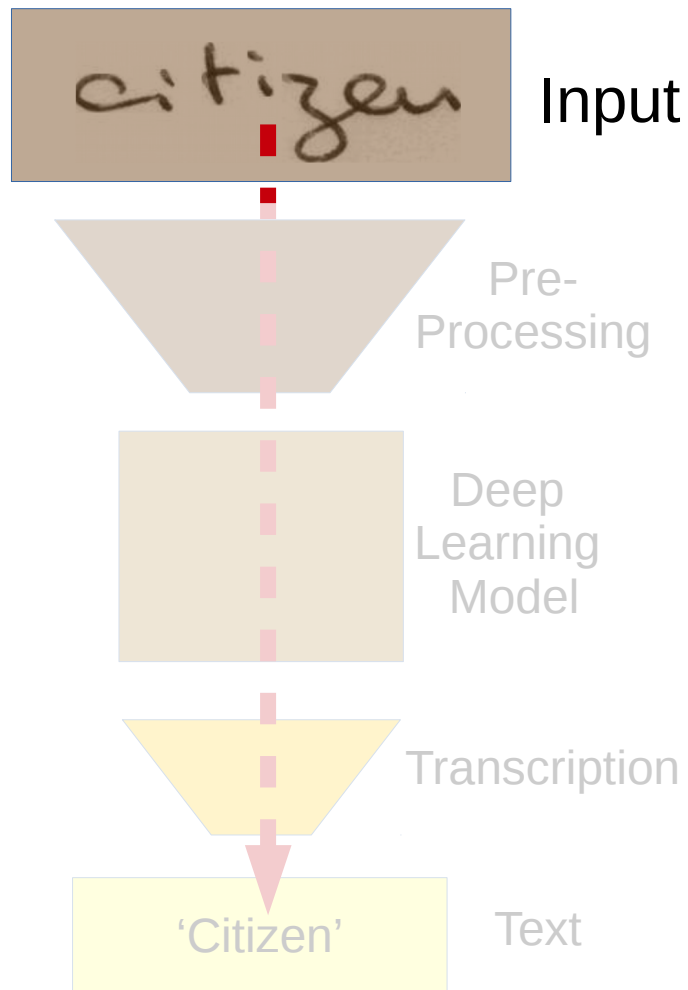# Goal: Transcription of Images to Text, without letter segmentation

- Examine a lot of different datasets, and found one which had the closest resemblance to my problem.

- Ended up with the IAM (offline, i.e. handwritten) dataset with **word segmentation**\*

- I'll skip over the first few months of research :)

- First let's see what the ML model is supposed to do:

\*The IAM dataset contains variations with word, sentence and line segmentation.

# Black Box Illustration



Input

Pre-Processing
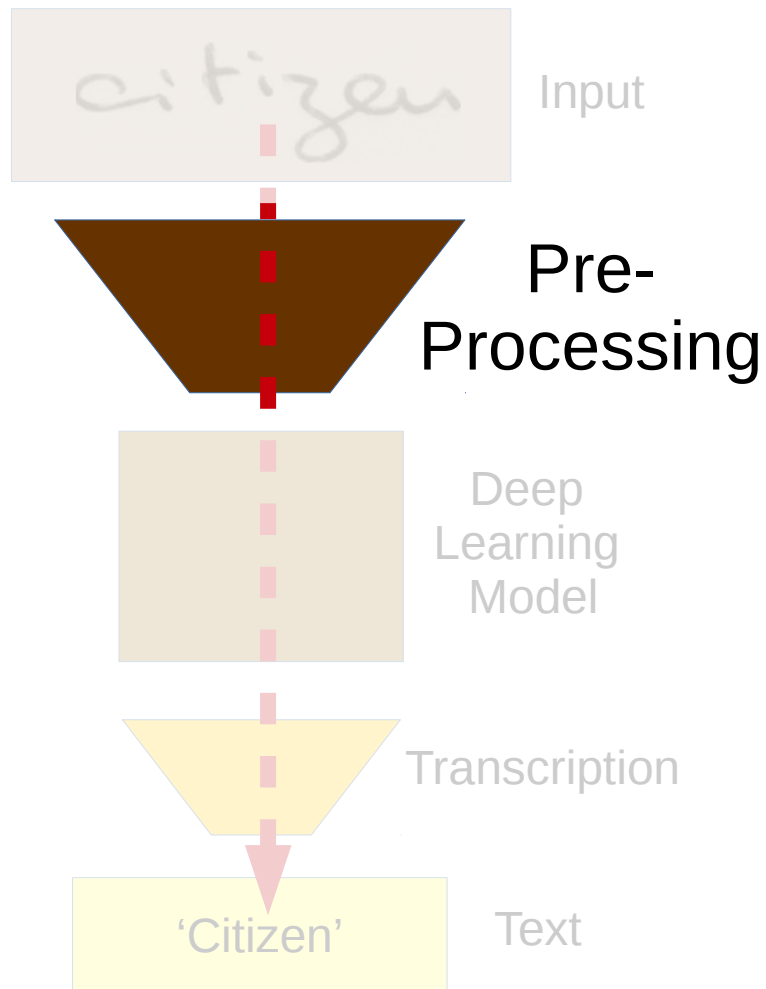
Deep Learning Model

Transcription

'Citizen'  Text

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'  Text

## Text Examples from IAM:



Labels from the Dataset:
'And' 'as' 'the' 'British'
'Government'
'stepped' 'up'

Images from the IAM dataset: http://www.fki.inf.unibe.ch/databases/iam-handwriting-database

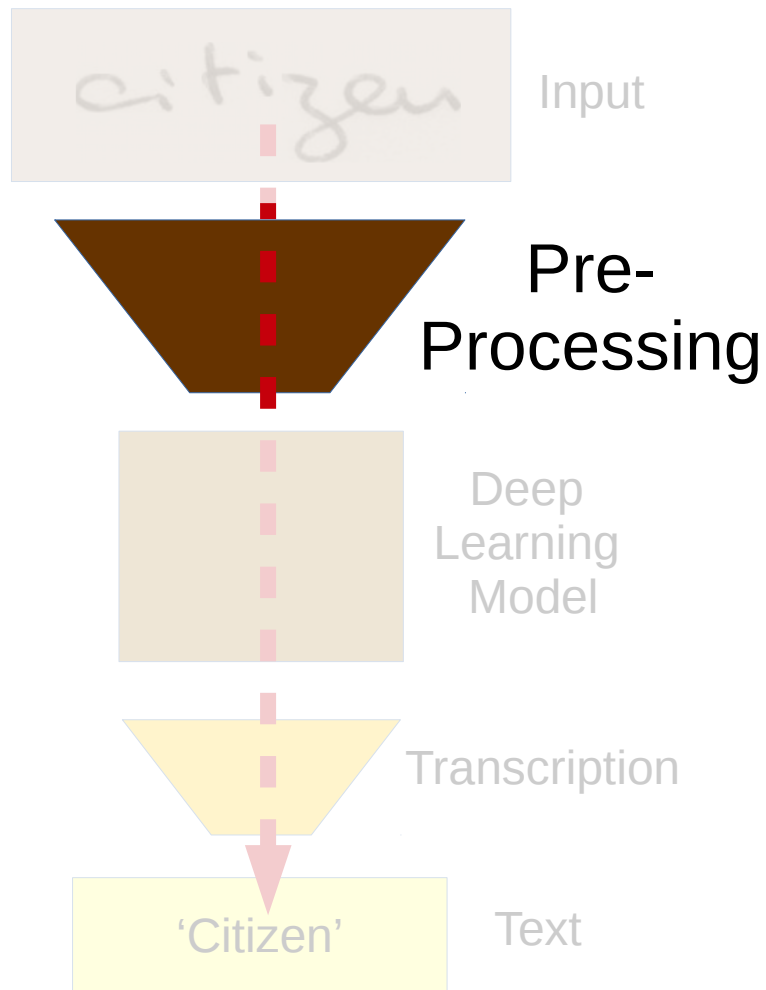# Black Box Illustration

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

Text had to be normalized to a fixed height:

Not ideal, but the letters should be recognized regardless of size anyway.

# Black Box Illustration

Input

Pre-Processing

Deep Learning Model
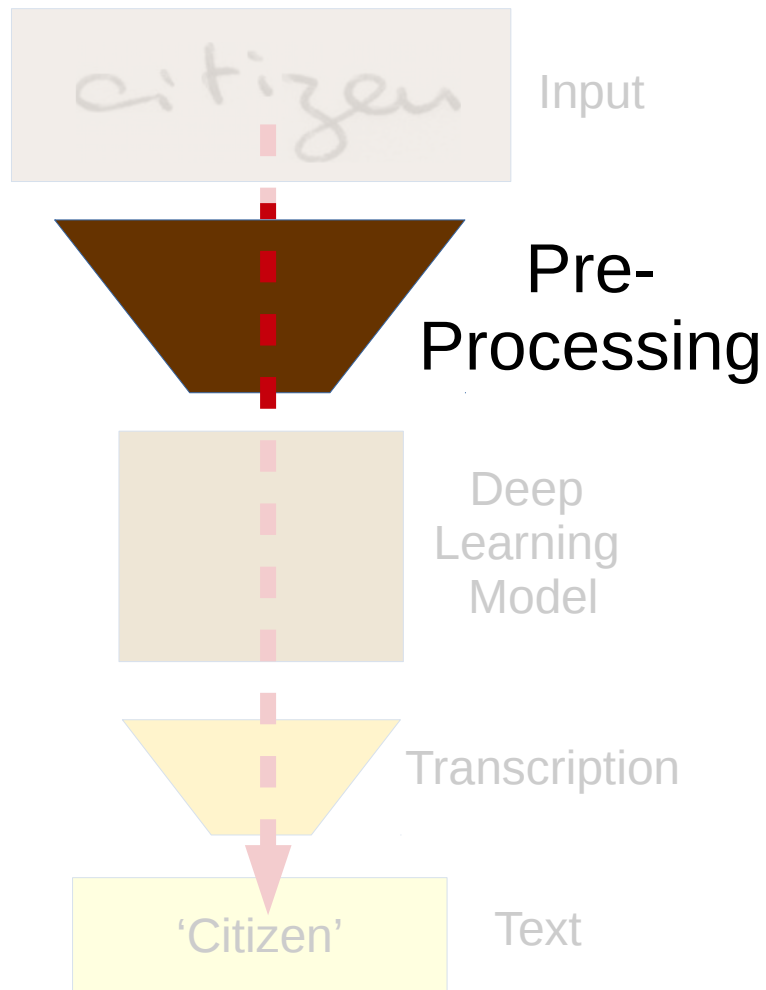
Transcription

'Citizen'

Text

Why not just pad the words?

Due to the large variation between the images padding might end up being very large, which would cause a visible break between the images and the padding. This is already a problem in some of the images:

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model

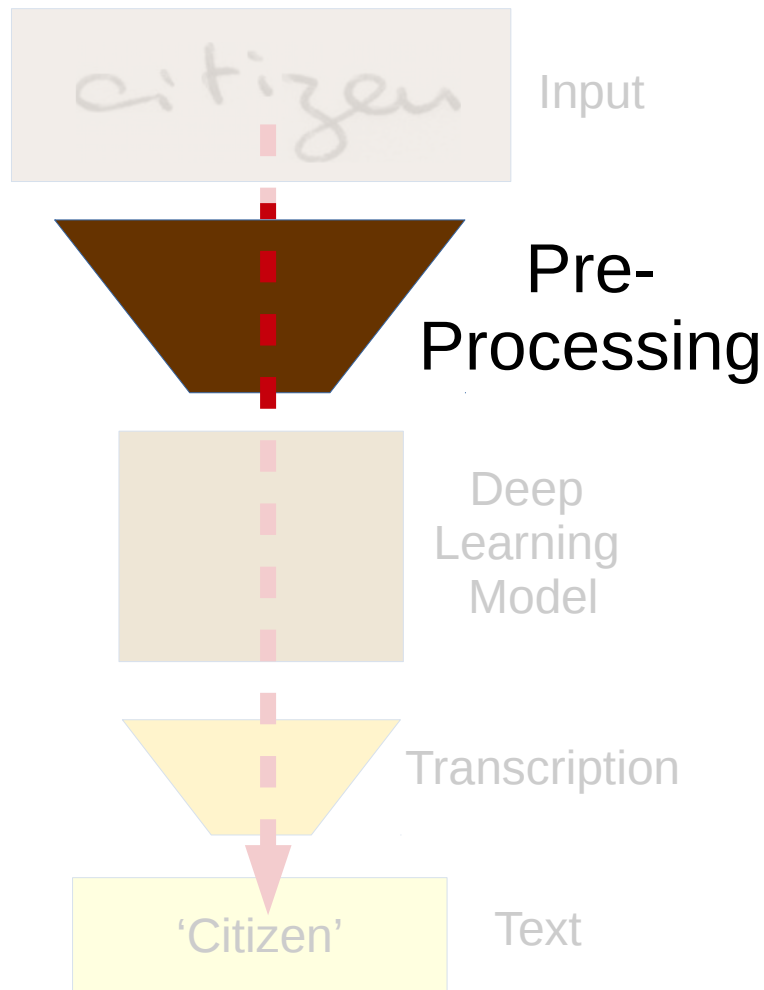Transcription

'Citizen'

Text



Different checks were made, the dataset was cleaned of the biggest errors in labels or images:

Some words were not segmented so that an entire sentence was labeled as one word, and sometimes the labels were wrong.
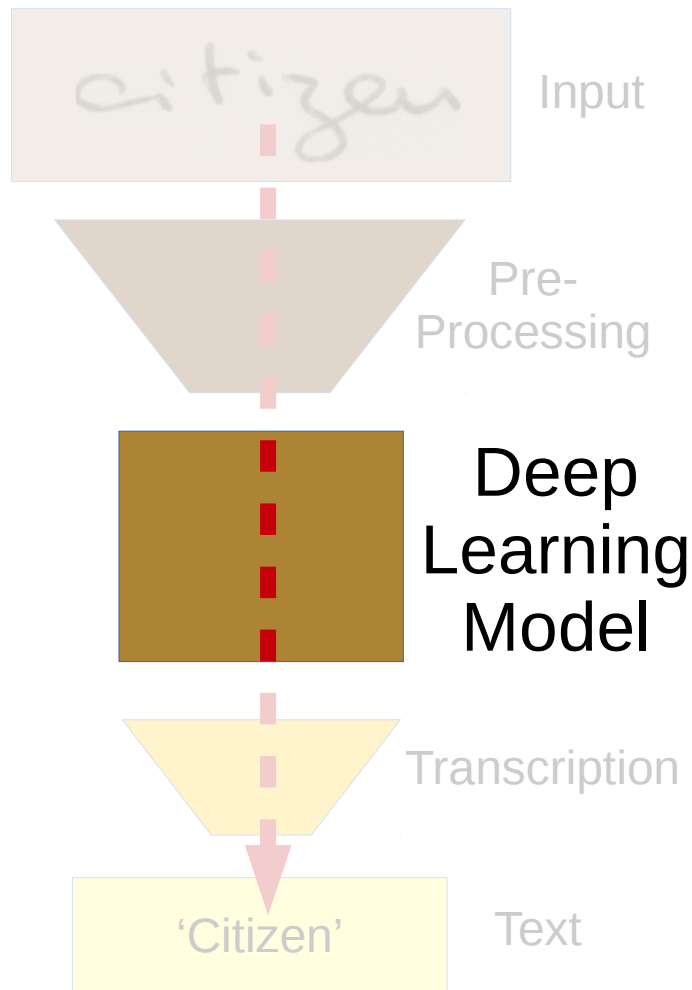


'And'

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

Images were pickled after preprocessing for easier and faster loading, as well as for the convenience of loading a single file instead of 115.276

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model
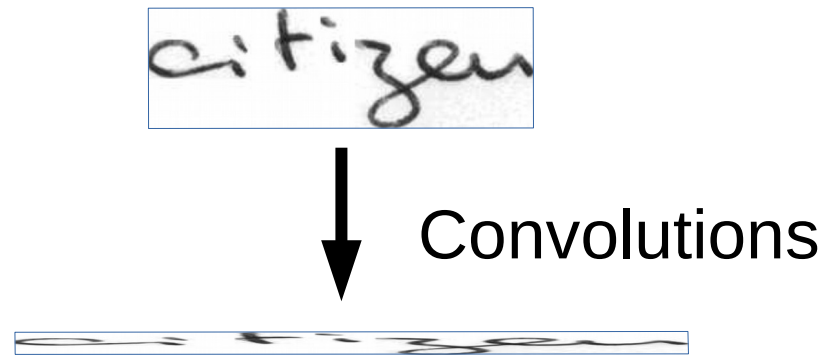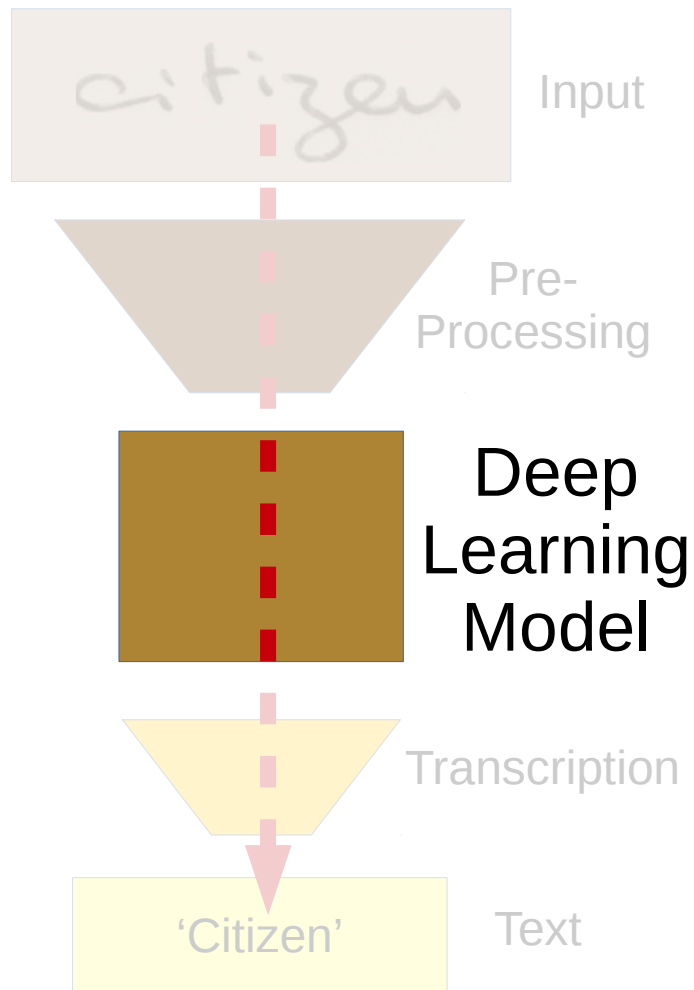
Transcription

'Citizen'

Text

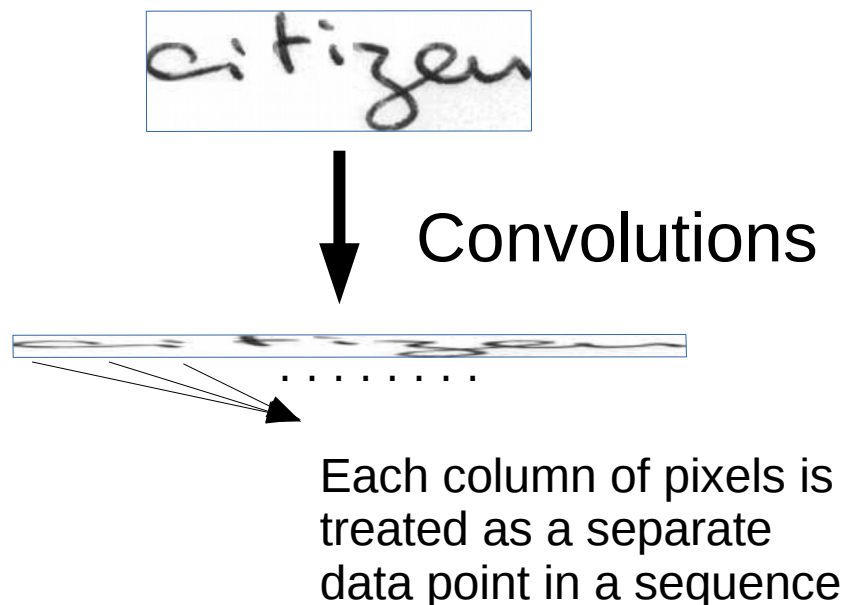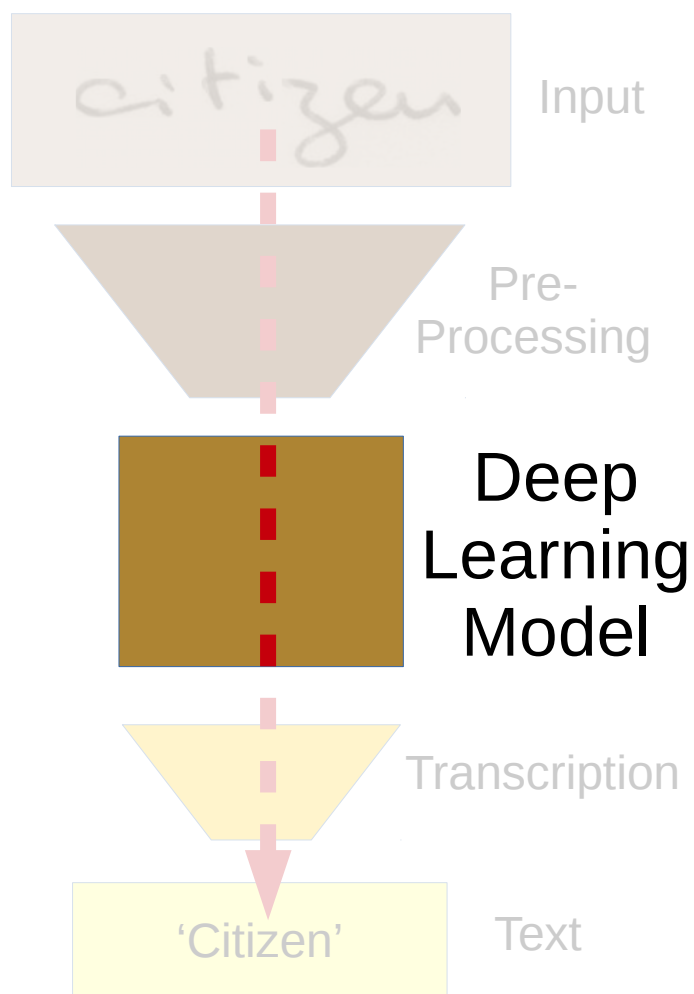After some research, I found this paper:
'An End-to-End Trainable NN for Image-based Sequence Recognition and it's application to Scene Text Recognition'.

What this paper, in a rather complicated fashion, describes, is  a model where Convolutional Layers are used to extract information from an image, and then a RNN is used to process that data to be transcribed into text.
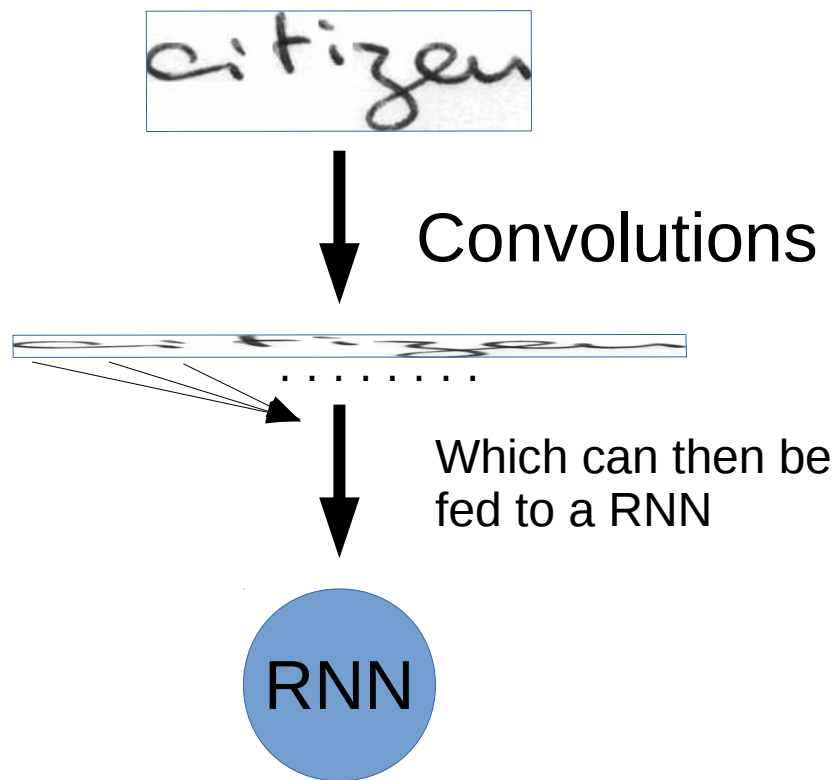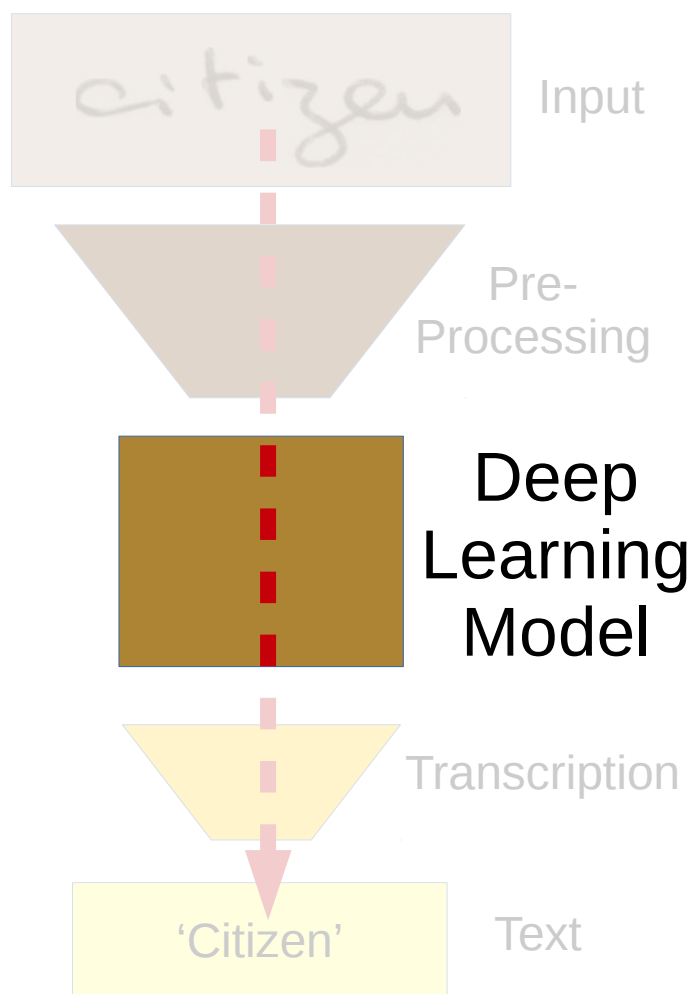
# Black Box Illustration

Input

Pre-Processing

## Deep Learning Model

Transcription

'Citizen'

Text

Convolutions

# Black Box Illustration

Input

Pre-Processing

**Deep Learning Model**

Transcription

'Citizen'

Text

Convolutions

. . . . . . . .

Each column of pixels is treated as a separate data point in a sequence

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

Convolutions

· · · · · · · ·

Which can then be fed to a RNN

RNN

# Black Box Illustration

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

Convolutions

· · · · · · · ·

Which can then be fed to a RNN

RNN

Transcription

# Black Box Illustration



Input

Pre-Processing

**Deep Learning Model**

Transcription

'Citizen'

Text

Convolutions

. . . . . . . .

Which can then be fed to a RNN

RNN

Transcription ➡ 'Citizen'

# Aside:

Input

Pre-
Processing

Deep
Learning
Model

Transcription

'Citizen'    Text

Time
dimension

Convolutions

. . . . . . . .

Which can then be
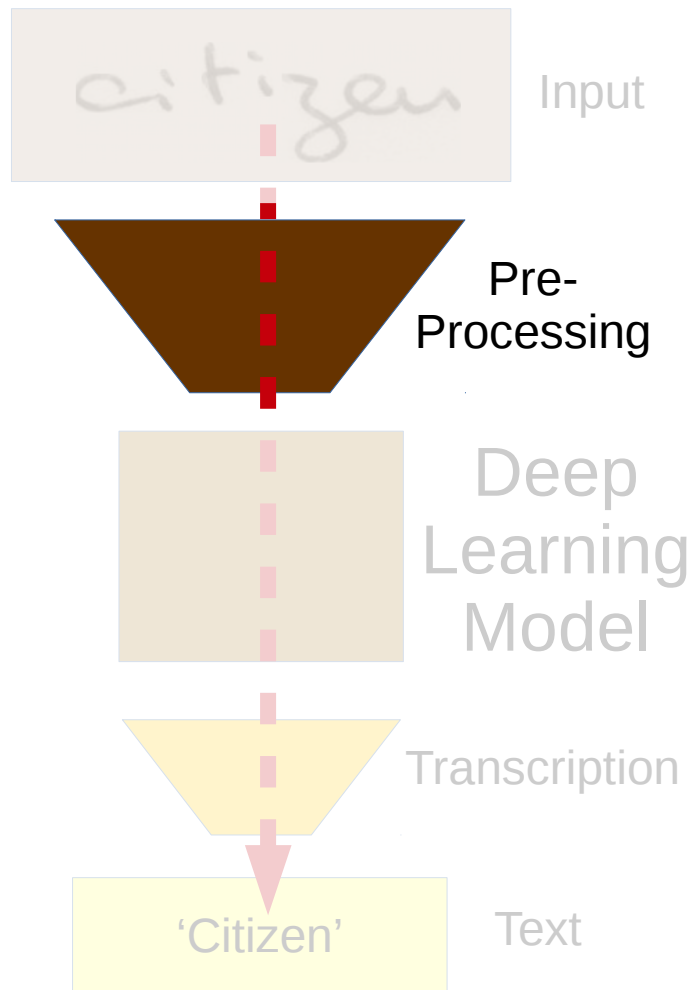fed to a RNN
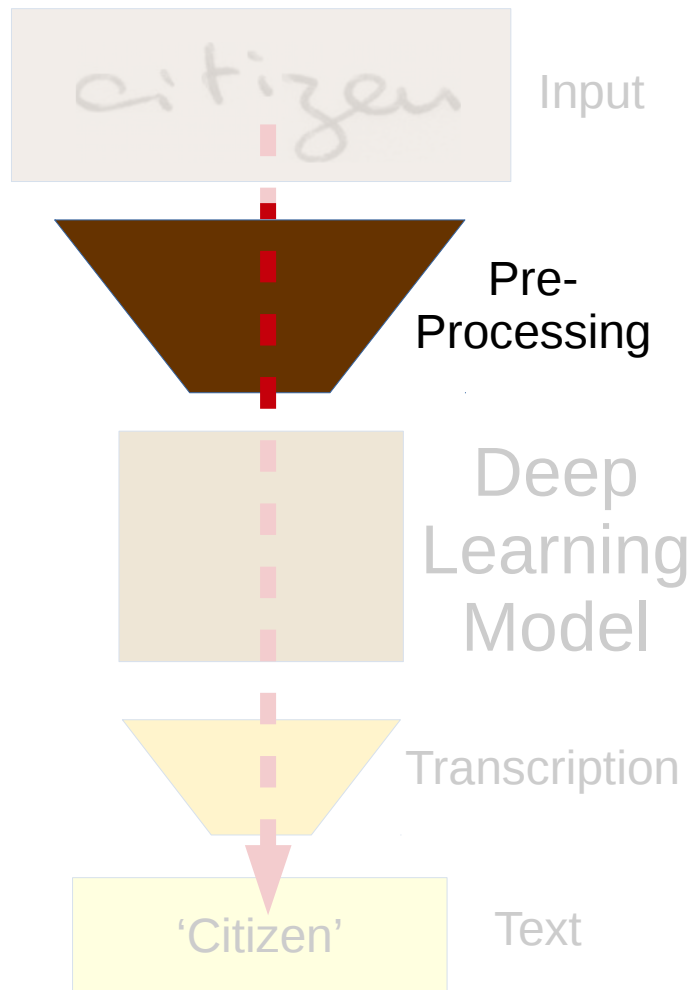
Since the RNNs take data
in which the 'Time
Dimension'* is the first
dimension, the
dimensions for the
images were flipped.

*See https://keras.io/layers/recurrent/

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

# Aside:

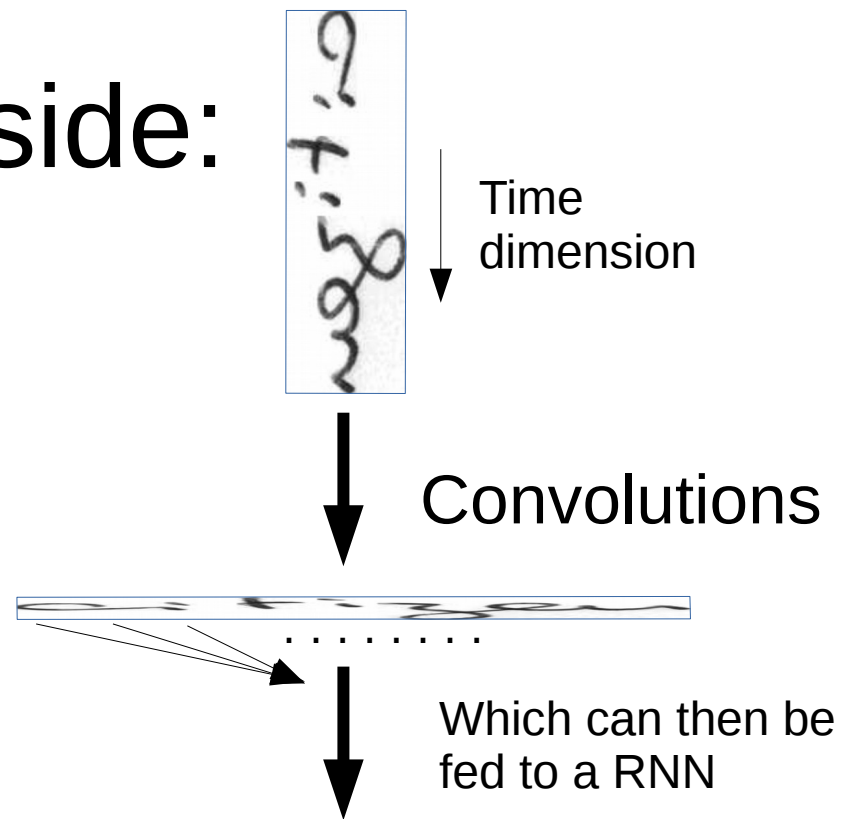Time dimension

Convolutions

. . . . . . . .

Which can then be fed to a RNN

Since the RNNs take data in which the 'Time Dimension'* is the first dimension, the dimensions for the images were flipped.
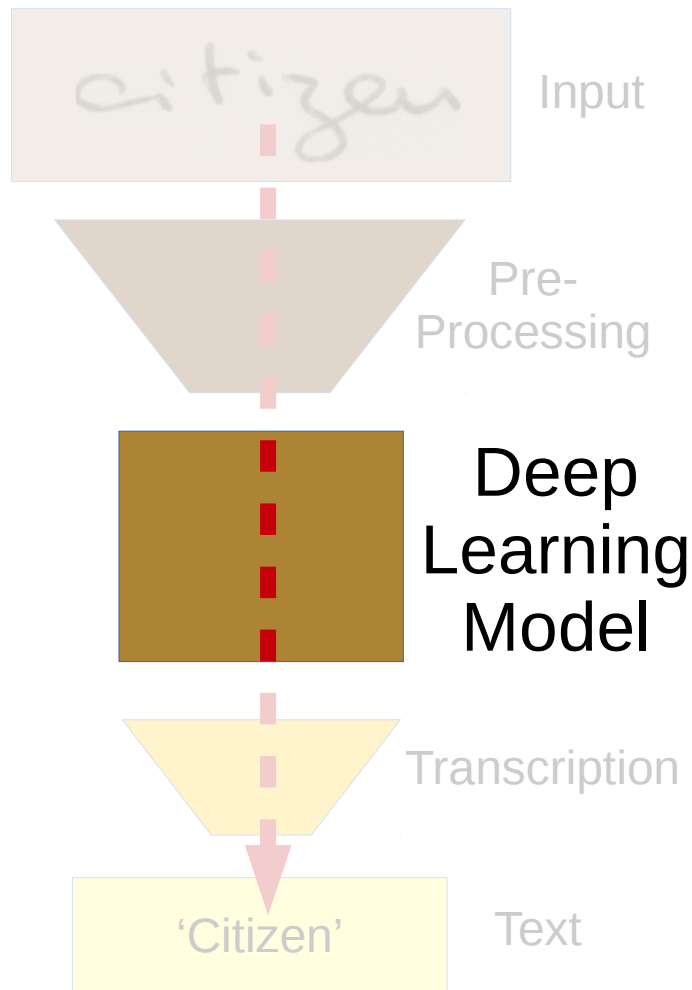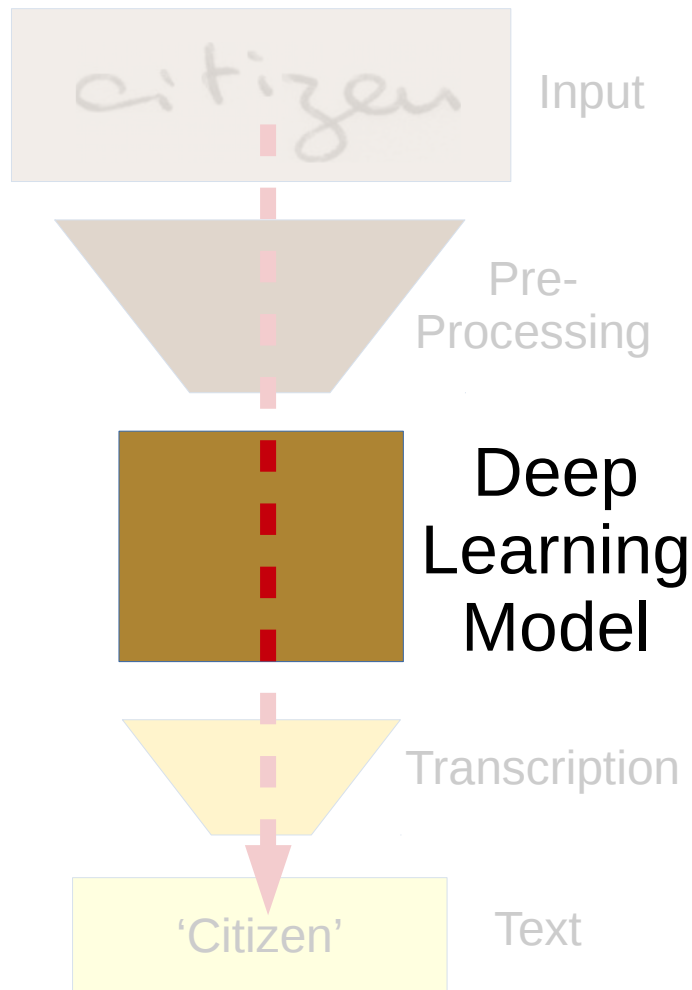
*See https://keras.io/layers/recurrent/

# Black Box Illustration

Input

Pre-Processing

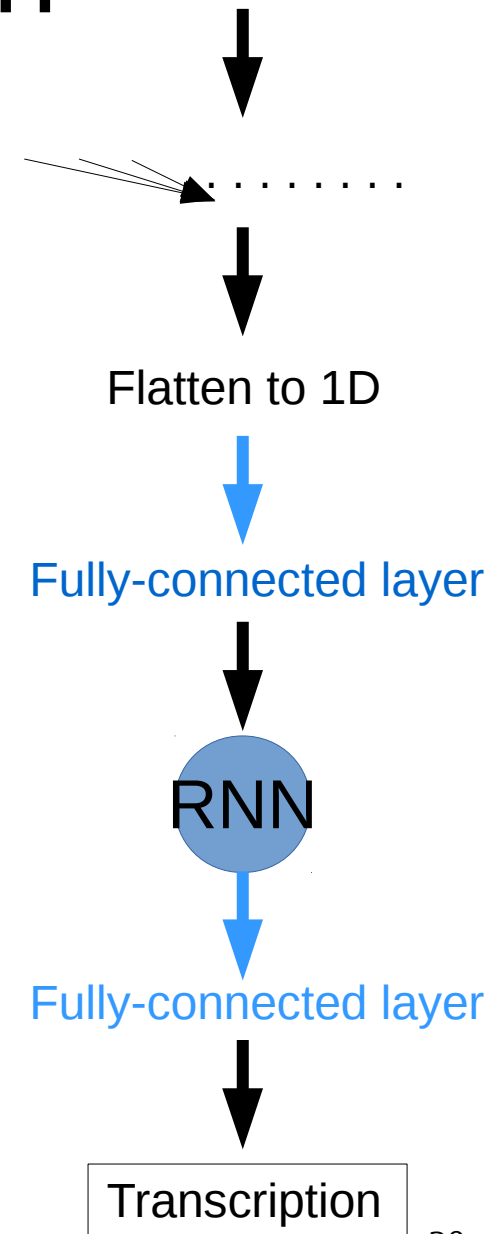Deep Learning Model

Transcription

'Citizen'

Text

Since the paper was coded in Torch I looked for any implementations that used Tensorflow and ideally the Keras frontend. While the architecture was relatively easy to understand, the puzzle piece missing was how to calculate the loss function for this type of output.
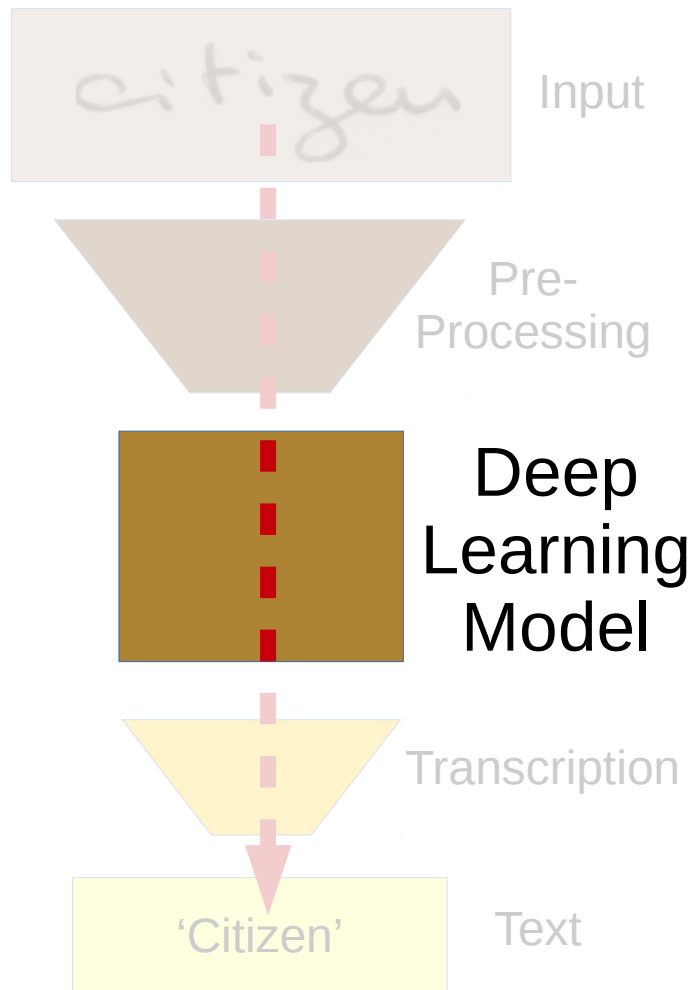
# Black Box Illustration

Input

Pre-Processing

## Deep Learning Model

Transcription

'Citizen'

Text

I found the keras-image-ocr repository* which (as a demonstration) used a similar architecture.

Flatten to 1D

Fully-connected layer

RNN

Fully-connected layer

Transcription

*https://github.com/Tony607/keras-image-ocr

# Black Box Illustration
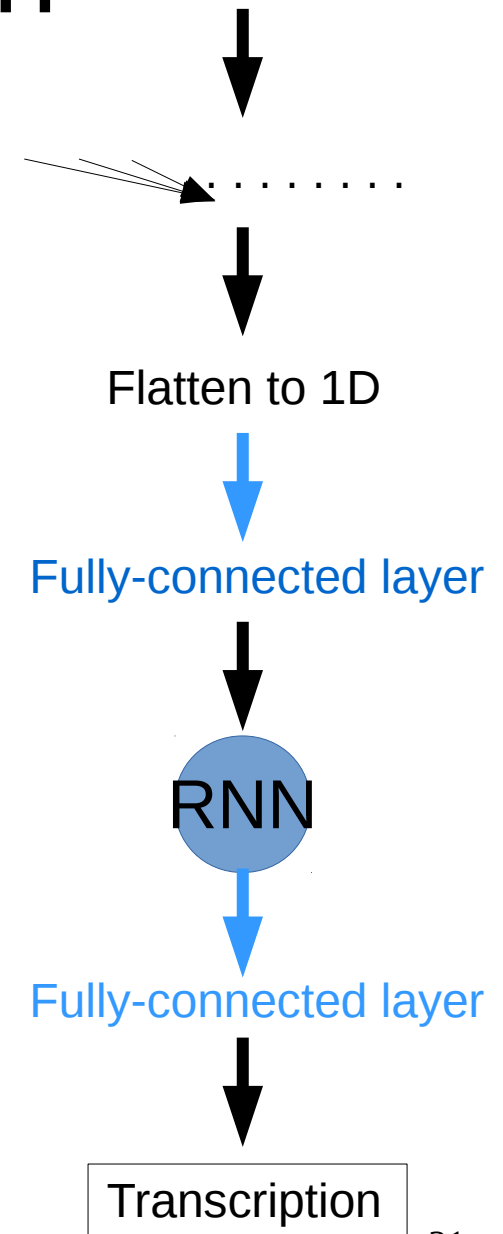
Input

Pre-Processing

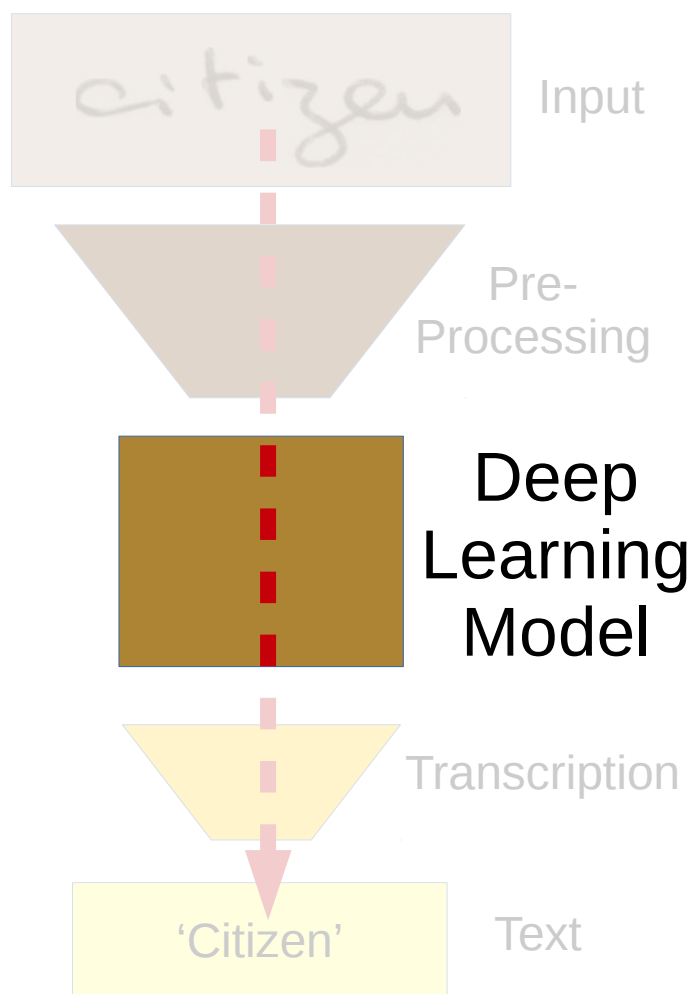Deep Learning Model

Transcription

'Citizen'

Text

I found the keras-image-ocr repository which (as a demonstration) used a similar architecture. This demo used Automatically generated images to train the network.
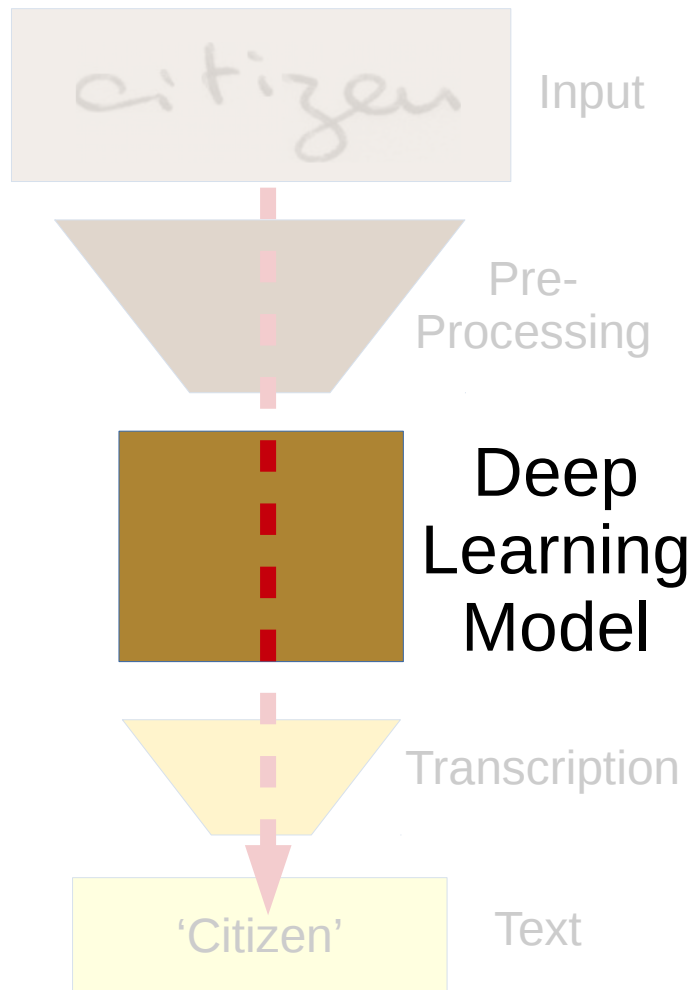
Example image*:

Flatten to 1D

Fully-connected layer

RNN

Fully-connected layer

Transcription

21

*https://github.com/Tony607/keras-image-ocr/blob/master/image-ocr.ipynb

# Black Box Illustration



Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

Most importantly it used a working implementation of the CTC loss function.

# Black Box Illustration
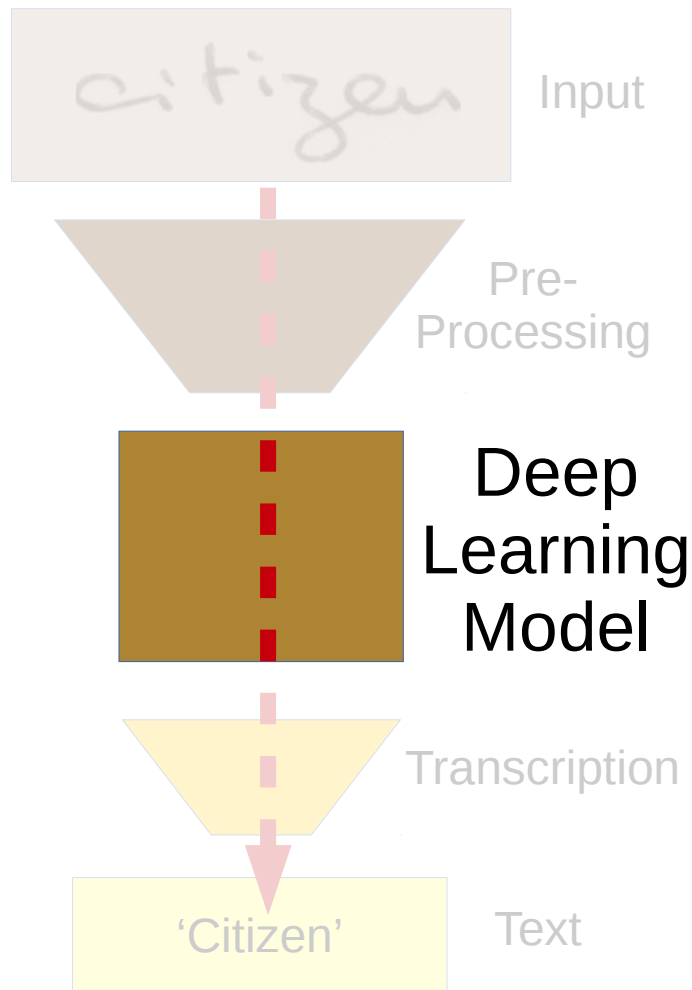
Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

Text

The CTC loss function is able to calculate Loss for unsynchronized input and output.
Since every letter is not the same number of pixels wide, you cannot apply a n-number of pixels per letter algorithm. This is where CTC comes into play. It was originally developed for audio data, but works for any kind of sequential data. I won't go into details as to how it works, see the paper at
http://www.cs.toronto.edu/~graves/icml_2006.pdf
for more. (warning: Math!)

# Black Box Illustration
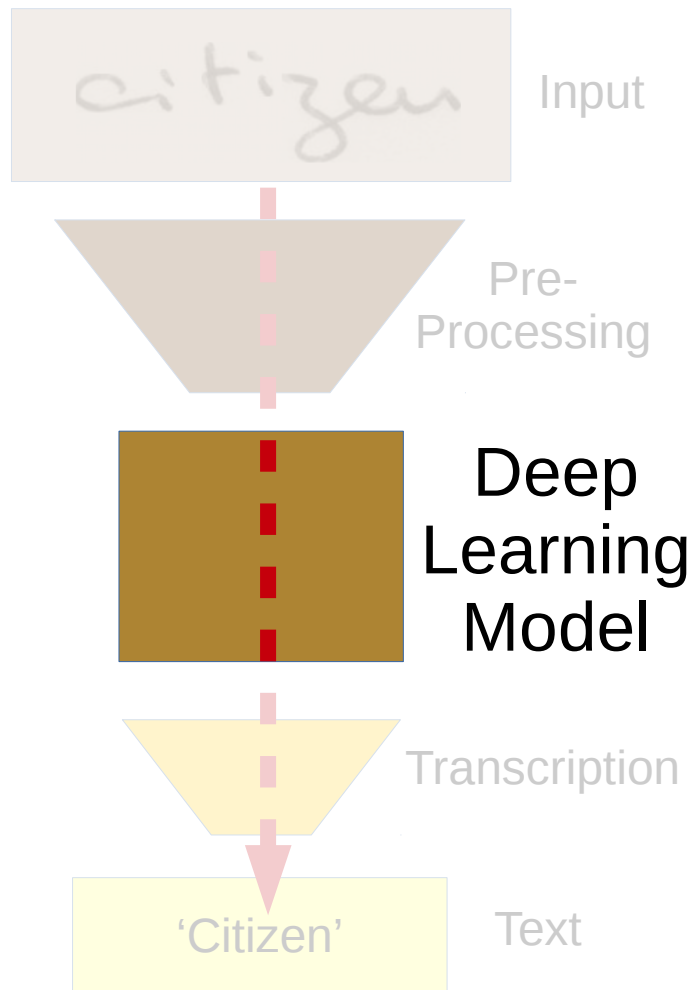
Input

Pre-Processing

Deep Learning Model

Transcription

Text

'Citizen'

While CTC Loss has been implemented in Tensorflow, the documentation is not particularly clear, so I had to take the data generator from the demo apart to see the data it was generating, then format the IAM dataset into that format.
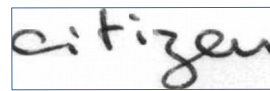
# Black Box Illustration

Input

Pre-Processing

Deep Learning Model

Transcription
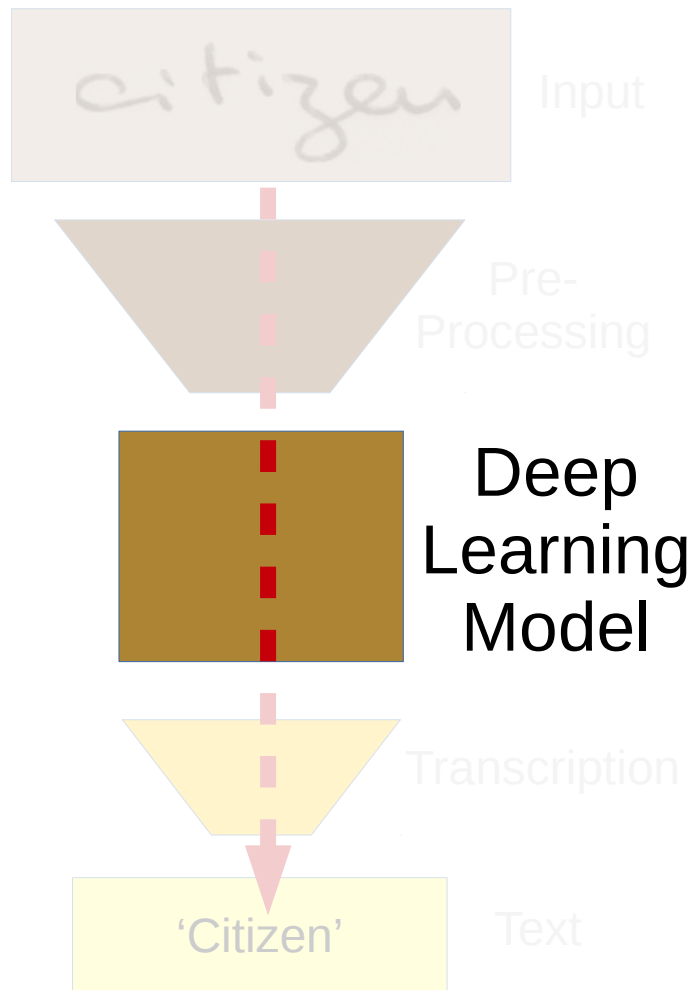
'Citizen'      Text

CTC loss requires for each sample:
- the output Data from the model,
- Labels encoded in integer form and padded with an encoding for 'blank' to a certain max length,
- The size of the input in the time dimension,
- The length of the label without the 'blank' padding.

For example:

(processed through the network),

'[5, 11, 22, 11, 28, 7, 16,-1,-1,-1]',

30 pixels,

7 letters long

# Black Box Illustration

Input

Pre-
Processing

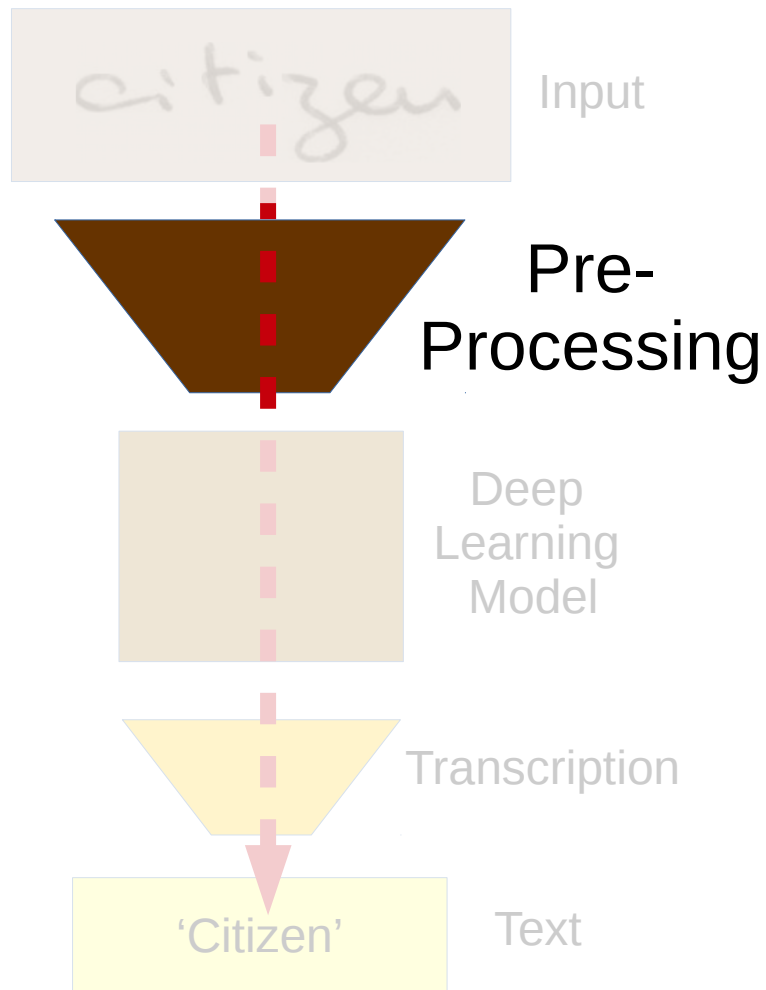Deep
Learning
Model

Transcription

'Citizen'     Text

Now the model can process images of variable width, however during training, image dimensions have to be defined, this is a restriction since gpu memory has to be reserved as well as computation optimizations needing a defined input format.
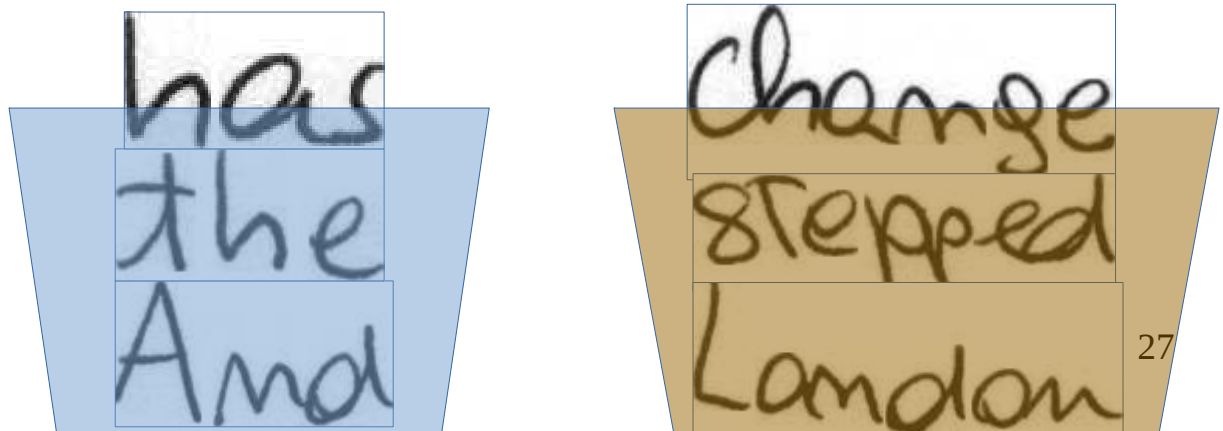
Solution: Padding?

Nope!

Back to data Pre-processing!

# Buckets!



Input

Pre-Processing

Deep Learning Model
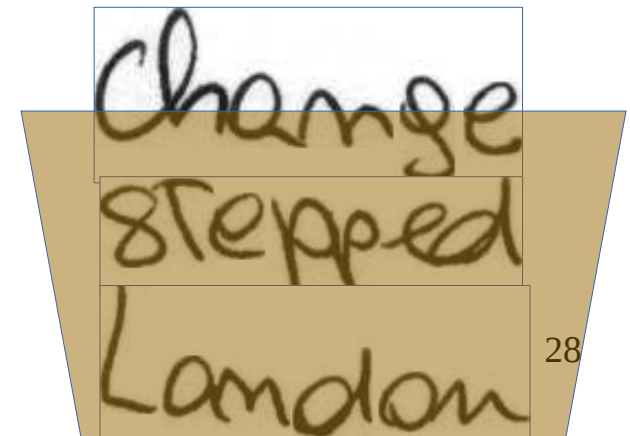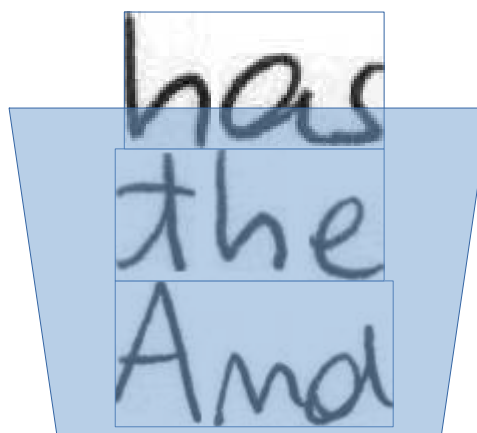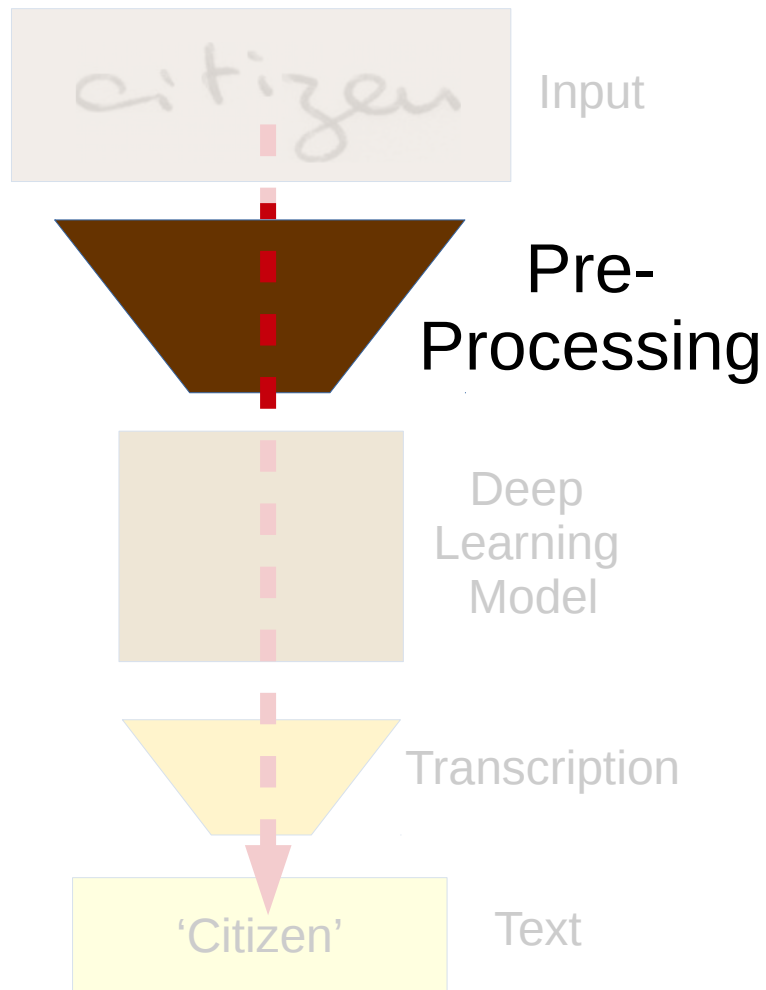
Transcription

'Citizen'

Text

Padding was not necessary since **Bucketing** was implemented. **Bucketing** is the grouping of images based on a criterion, in this case **Image Length**. Image height was already normalized as explained before.

# Buckets!

Input

Pre-Processing

Deep Learning Model
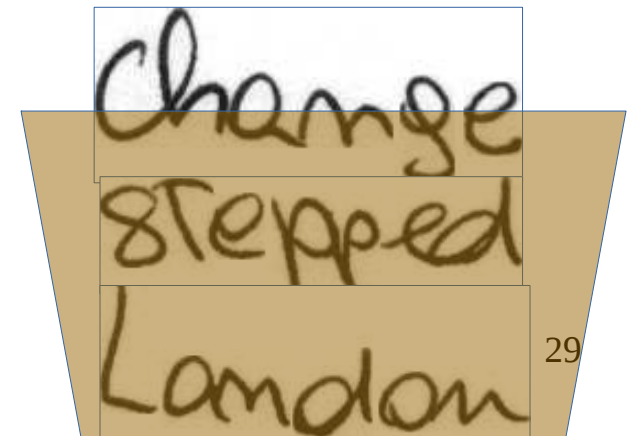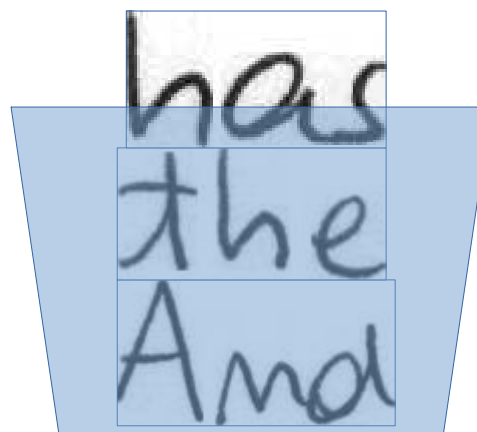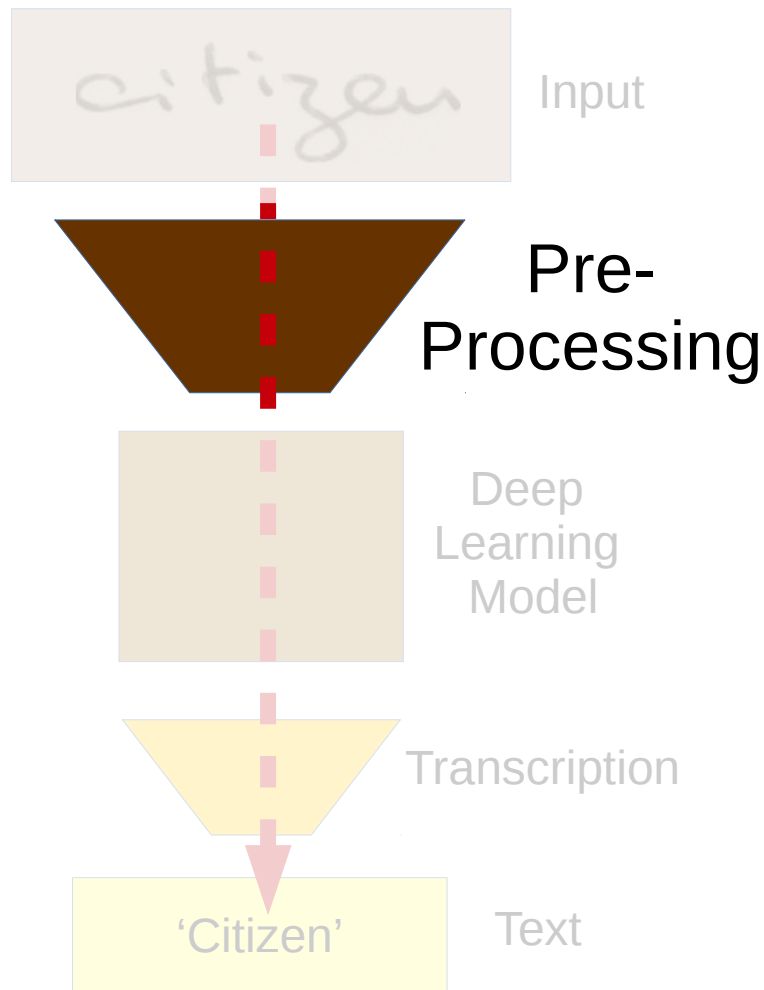
Transcription

'Citizen'                Text

Now we can feed a bucket at a time to the model and since all the images have the same width, the model can be trained on a bucket at a time.
All that is left is writing a generator which supplies the model with a bucket at a time. This can probably be improved with Queues, but it works well enough for now.

# Buckets!

Input

Pre-Processing

Deep
Learning
Model

Transcription

'Citizen'

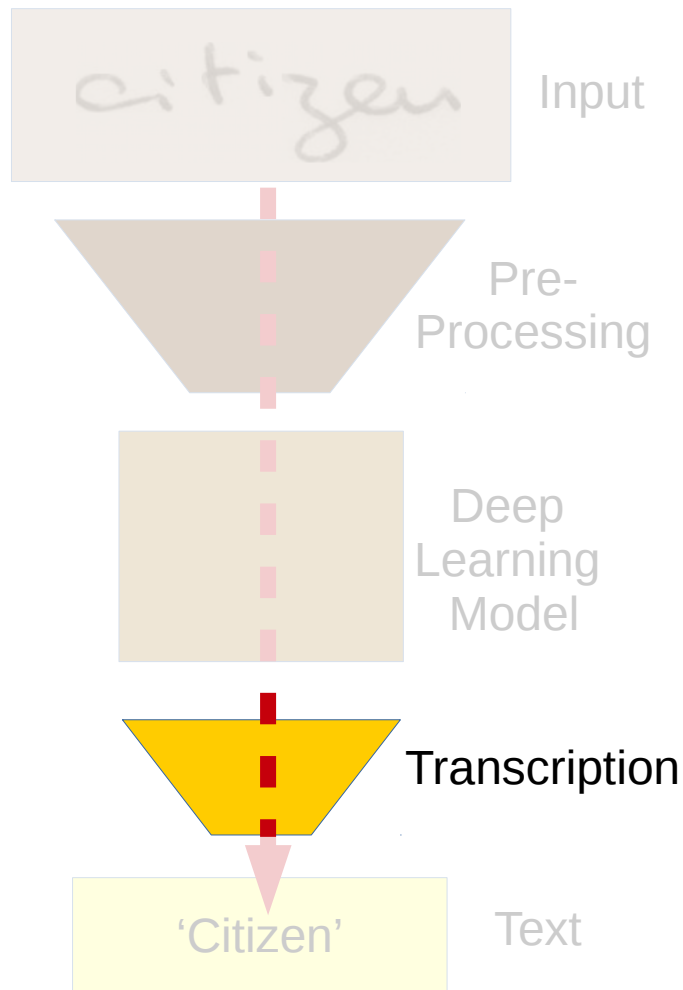Text

One last thing about the generator: Keras runs the generator in the background on the CPU, while the Model is training. This saves time by pre-fetching data. So even if the generator only supplies a single bucket,

# Black Box Illustration

Input

Pre-Processing

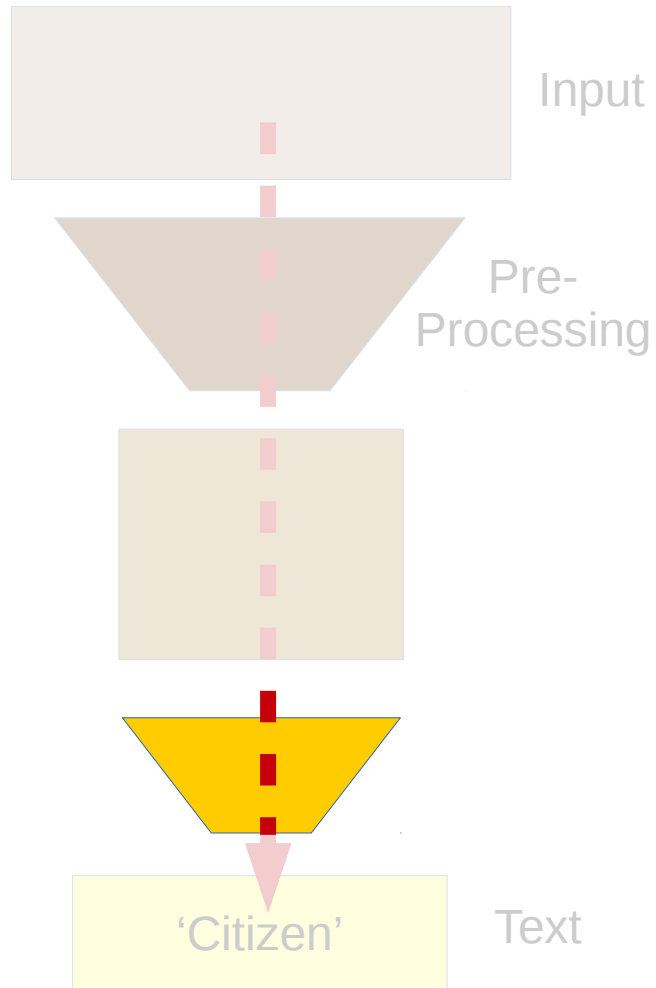Deep Learning Model

Transcription

'Citizen'   Text

Transcription is done via a CTC decode function from the Keras frontend, which uses Beam Search.

Beam search is a transcription method where the most probable paths are explored and the highest scoring end path is used.

30

# Beam Search, top_paths = 3

Input

Pre-
Processing

'Citizen'    Text

Example: 'The dog walked'

'A'    'The'    'dog'

# Beam Search, top_paths = 3

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

Example: 'The dog walked'

'A'    'The'    'dog'

'Dog' 'Tree' 'Man'

'Dog' 'Tree' 'Man'

'walks' 'runs' ' eats'

# Beam Search, top_paths = 3

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

Example: 'The dog walked'

'A'        'The'        'dog'

'Dog'  'Tree'  'Man'
P= 0.6    0.2    0.4

'Dog' 'Tree' 'Man'
P= 0.7    0.2    0.5

'walks' 'runs' ' eats'
P= 0.2    0.1    0.2

# Beam Search, top_paths = 3

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

Example: 'The dog walked'

'A'      'The'      'dog'

'Dog' 'Tree' 'Man'        'walks' 'runs' ' eats'

'Dog' 'Tree' 'Man'

'A Dog'        'The  Dog'        'The  Man'
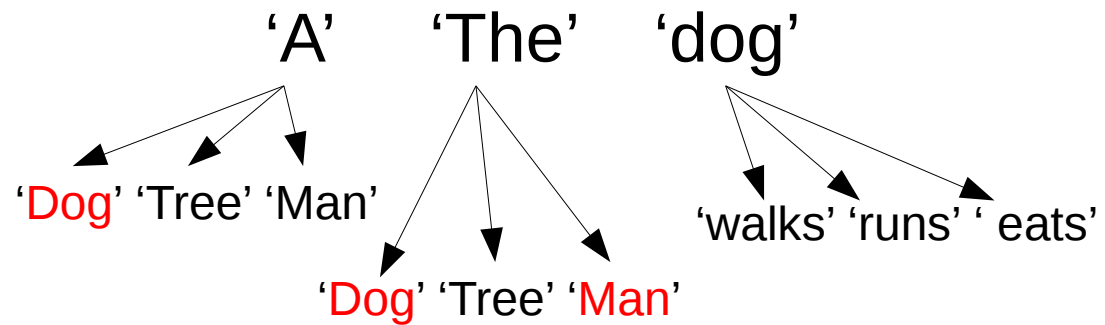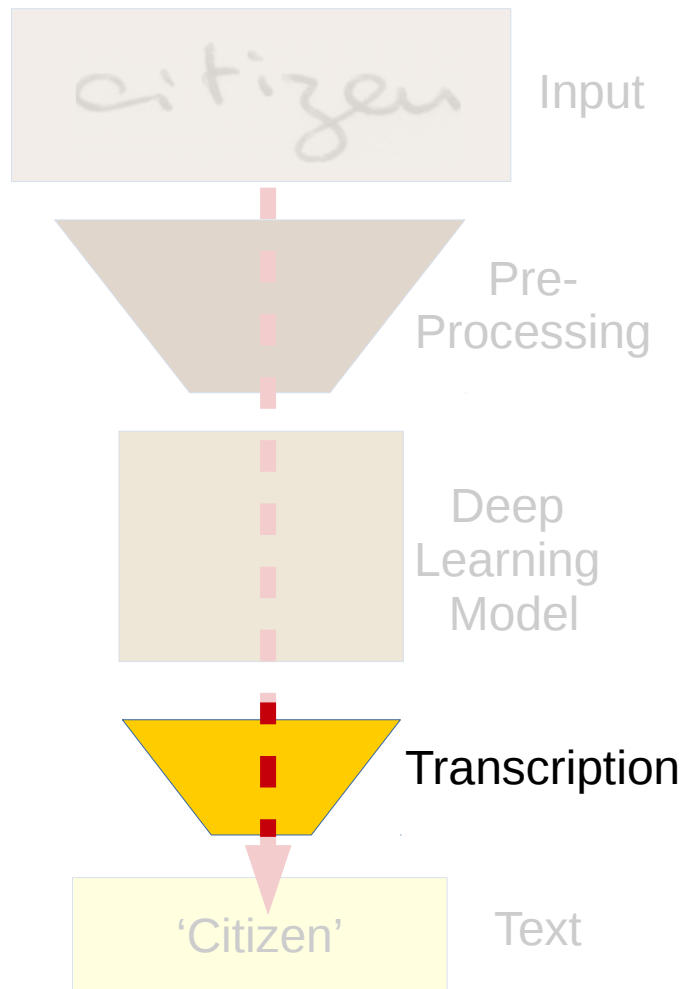
# Beam Search, top_paths = 3



Example: 'The dog walked'

'A'    'The'    'dog'

'Dog' 'Tree' 'Man'        'walks' 'runs' ' eats'

'Dog' 'Tree' 'Man'

'A Dog'        'The  Dog'        'The  Man'

'walked'  '<END>' 'climbed'
P= 0.7   0.1   0.3

'walked'  'talked' '<END>'
P= 0.8   0.1   0.3

'walked'  'talked' 'climbed'
P= 0.5   0.6   0.4

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

# Beam Search, top_paths = 3

Example: 'The dog walked'

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'

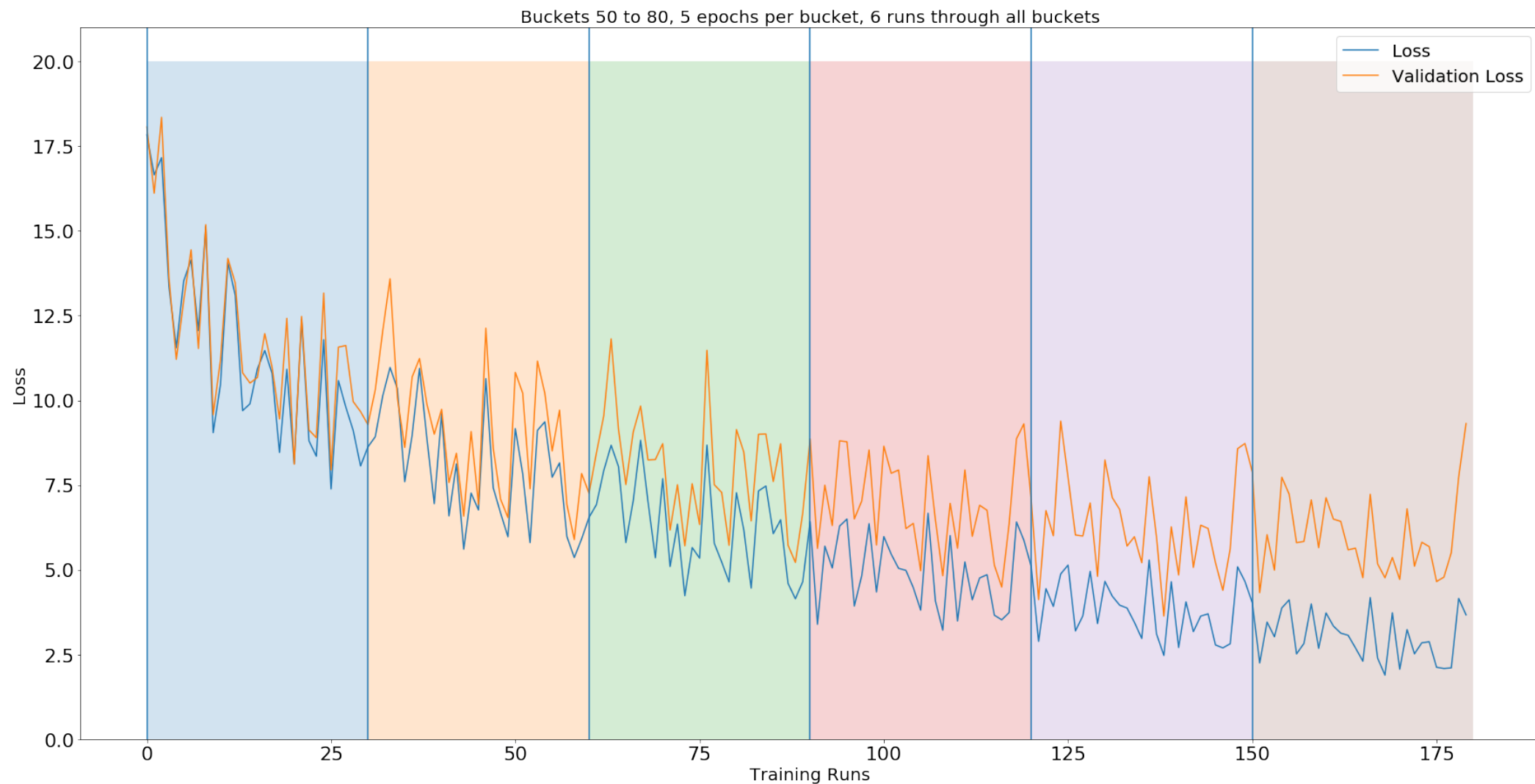'A'   'The'   'dog'

'Dog' 'Tree' 'Man'         'walks' 'runs' ' eats'

'Dog' 'Tree' 'Man'

'A Dog'        'The  Dog'        'The  Man'

'walked'  '<END>' 'climbed'        'walked'  'talked' '<END>'        'walked' 'talked' 'climbed'

'away'  <END>  'towards'        'to'  <END>  'with'

'away'  <END>  'towards'        P = 0.6

....        P = 0.8

# Training



Buckets 50 to 80, 5 epochs per bucket, 6 runs through all buckets
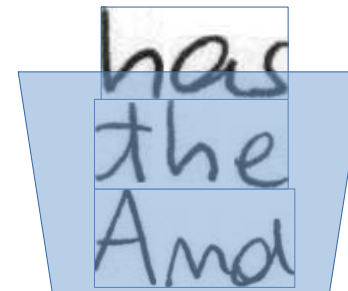
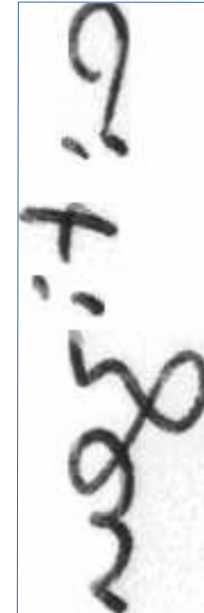Small bucket number causes the model to quickly begin to overfit

# Training



Buckets 100 to 200, 2 epochs per bucket, 4 runs through all buckets
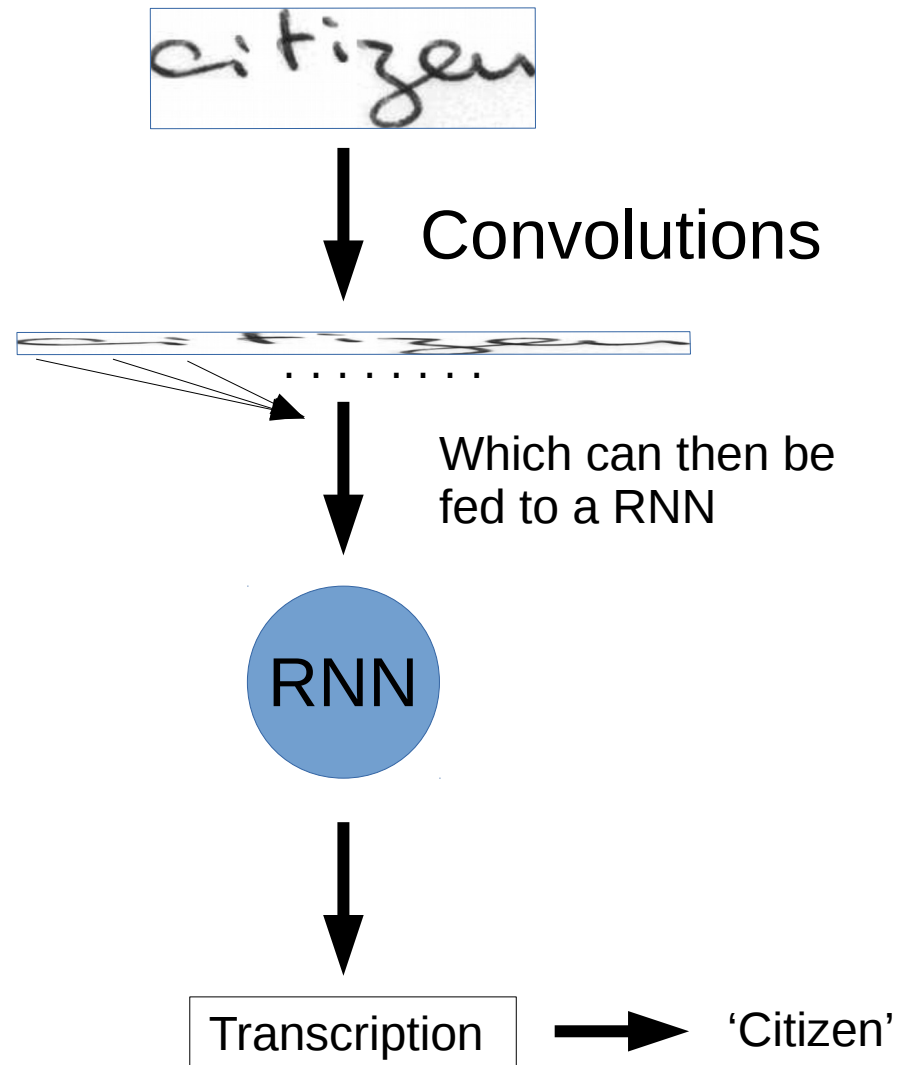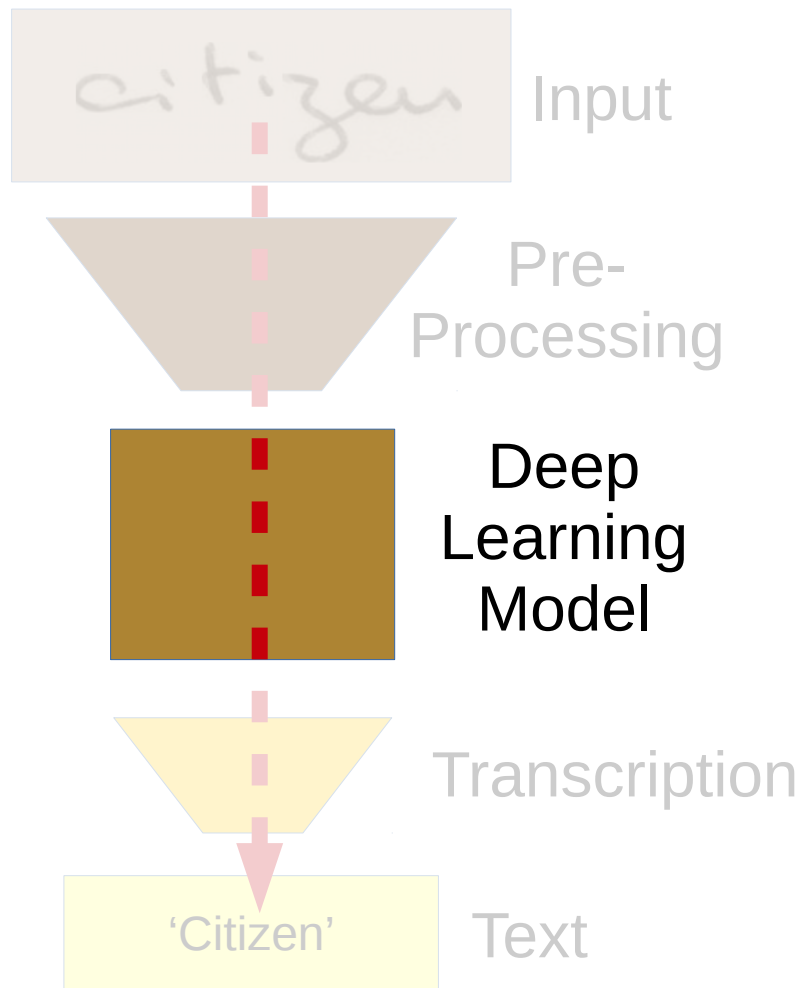
Only 70 of 100 buckets were used, since 30 did not contain enough data for validation. Larger diversity of data and fewer epochs combat overfitting.

38

# QUESTIONS ?



Input

Pre-Processing
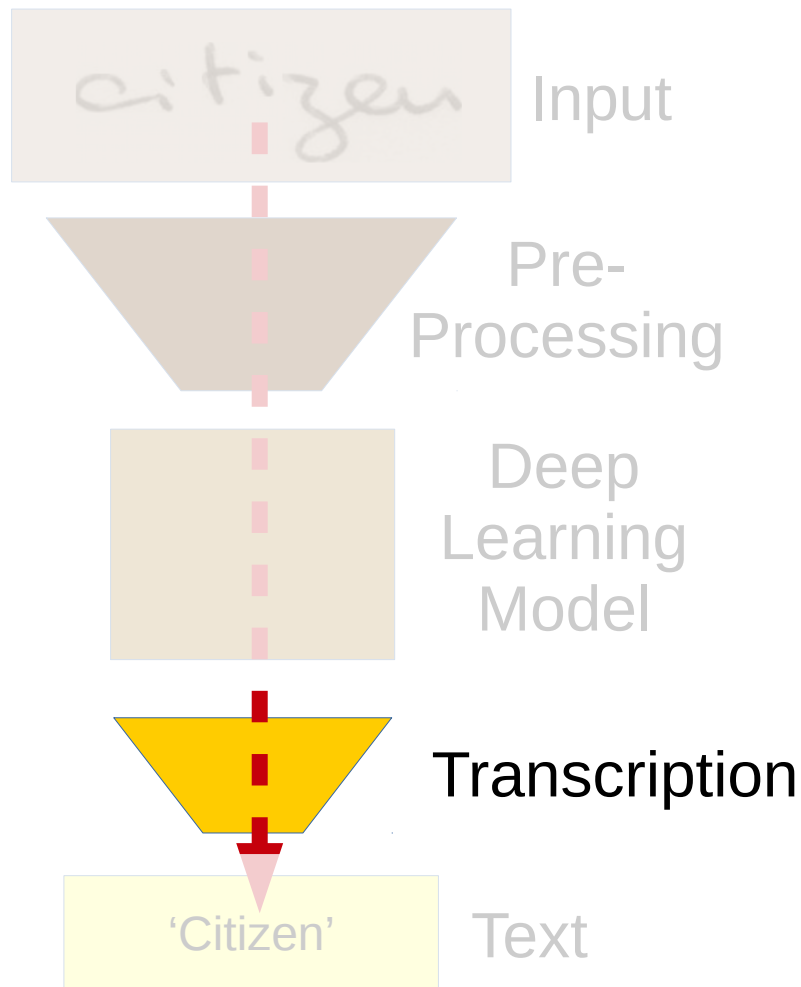
Deep Learning Model

Transcription

'Citizen'

Text

# QUESTIONS ?



Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'     Text

Convolutions

. . . . . . . .

Which can then be fed to a RNN

RNN

Transcription   →   'Citizen'

# QUESTIONS ?

Input

Pre-Processing

Deep Learning Model

Transcription

'Citizen'    Text

## Beam Search

'A'    'The'    'dog'

'Dog' 'Tree' 'Man'

'The  Dog'

'walked'  'talked' '<END>'

'away'  <END>  'towards'

P = 0.8