

# Deep-Learning Architectures for Handwritten Text Recognition

## Concepts & Building Blocks

Pablo M. Olmos  
olmos@tsc.uc3m.es



# Index

- 1 Motivation
- 2 Recurrent Neural Networks for Supervised Sequence Learning
- 3 Connectionist Temporal Classification
- 4 A DL network for Handwritten Text Recognition

si Re v

13

9

Capit<sup>on</sup> con Juan ponce de leon  
Sobre el descubim<sup>to</sup> de la ysla  
Beniny —

Por quant<sup>o</sup> Vos Juan ponce de leon me embiastes a suplicar e pedir por merced  
vos diese licencia y facultad para ir a descubrir y poblar la ysla de Beniny  
con ciertas condiciones que adelante sean declaradas. Por ende por vos  
hazer merced vos doy licencia y facultad para que podais ir a descubrir y  
poblar la dicha ysla con tanto que no sea de las que hasta agora estan  
descubiertas. y con las condiciones y segun que adelante sea contenido.  
En esta guisa

CAPITULACION CON JUAN PONCE DE LEON SOBRE EL DESCUBRIMIENTO DE  
LA ISLA DE BENINY.

AÑO DE 1512 (1).

Por quanto vos Juan Ponce de Leon, Me embiastes a suplicar e pedir por merced, vos diese licencia y facultad para ir a descubrir y poblar las Islas de Beniny, con ciertas condiciones que adelante seran declaradas, por ende, por vos hacer merced, vos doy licencia y facultad para que podais ir á descubrir y poblar la dicha Isla, con tanto que no sea de las que hasta ahora estan descubiertas, y con las condiciones y segun que adelante será contenido en esta guisa.

confirmación

29  
25  
2  
e Rey

Por quanto Vos el conde don Hernando de andiada y Xpo Valdeharo Nro factor de la casa de la contratación de la especiería me hicisteis Relación que por Nos servir quereis hacer cierto viaje y descubrimiento. En las nuestras Indias del mar oceano dentro de los límites y tierras de demarcación y que para ello armariades con las condiciones que de yuso se han contenidas. Una carabela de porte de cinquenta hasta sesenta toneles. y un patax de Veinte y cinco a treinta toneles. forneados de las cosas necesarias. Asideca por vos como de mar y mantenimiento y otras cosas que se requieren para semejante viaje y descubrimiento y que de mas de la dicha carabela y patax. enviareis En pieças Vn bergantín de naves para descubrir qualquier riberia por las partes don aveare. y me suplicasteis y pedistes para mrd Vos mandásemos. de licencia y facultad para ello. E yo por vos hacer mrd tobo lo por bien y sobre ello mande tomar con vos otros. La sientro y capitulación siguiente.

CAPITULACION QUE SE TOMÓ CON HERNANDO DE ANDIADA Y CRISTÓBAL DE HARO PARA HACER VARIOS DESCUBRIMIENTOS.

AÑO DE 1526 (1).

Por quanto vos el Conde Don Hernando de Andia y Cristóbal de Haro, Nuestro factor de la casa de la contratación de la especiería, Me hicisteis relación, que por Nos servir, quereis hacer cierto viaje y descubrimiento en las Nuestras Yndias del mar Oceano, dentro de los límites y tierras de Nuestra demarcación; y que para ello, armariades con las condiciones que de suso serán contenidas, una carabela de porte de cinquenta hasta sesenta toneles y un patax de veinte y cinco á treinta toneles, forneidos de las co-



La Reyna

r sero seris de ob poe ob affo de finto lren  
saber a miced obio tulo congo seado  
pangoloea amlada de la epe e y ad ee  
Cnoba affo y dm de e per ad r m omu ge  
m do finto am gerd ob e g d o y m  
m d n l a e e h e a m m d o n p e o y o n d  
m d o r t u m o r t a d e r e n e n t a b r m e s e d i  
o t r a o r d e s a e m p a m o s d e e a e y n d i a s  
d e n g o d e e s t m i p e n a a m m p o s e r  
o d i e e e a s d i a o l m e n t a s d e m i s e l a d i a  
o d n e a o r e a m m o p e r o b o e p h e m t y f r a d o  
m a v e s b a y d o m i e t r a d o a d o n l a s d e n o f n e n  
t a s a m m o p e r a d e s t a m m m m e d i a o m i  
o m p a s s a d e p o s e e r v o l u e m a n d o s e m p  
d e b e m e d i a s p i s d i x i m o d e e p n e s o f f a

There are 80 million of documents at the Archivo General de Indias de Sevilla. Only the 10% are digitalized and, **among them**, only the 10% of them are transcribed!

Researchers have to rely on the metadata provided for each document by the archivists. It contains very general information about the topic of the document.

### Our (very) long-term goals

- Transcription of the millions of historic documents stored at historic archives.
- Provide search-by-content tools.
- Find latent structures in the corpus of documents, able to find related documents beyond what the current document metadata might suggest.

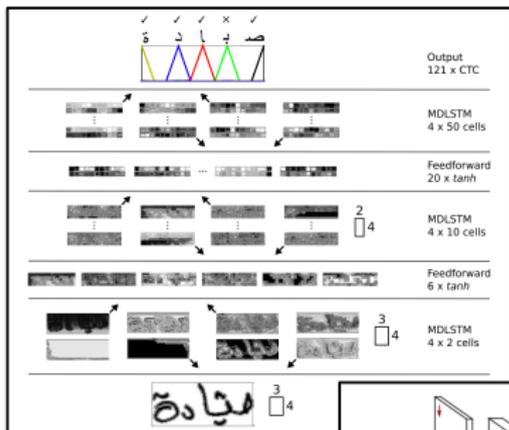
## Current Status of the Project

- We are now implementing state-of-the-art techniques based on Deep-Learning (DL) for Modern Handwritten Text Recognition.
- **IAM Database** of (labelled) handwritten text.
- Investigate generalization to historical documents:
  - ▶ **Very different handwriting styles.** Almost each document has been written by a different person ...
  - ▶ Only a few pages of each document.
  - ▶ How much labeled data are enough?
  - ▶ ...

# Today

Describe the basic elements of DL architectures for Handwritten Text Recognition:

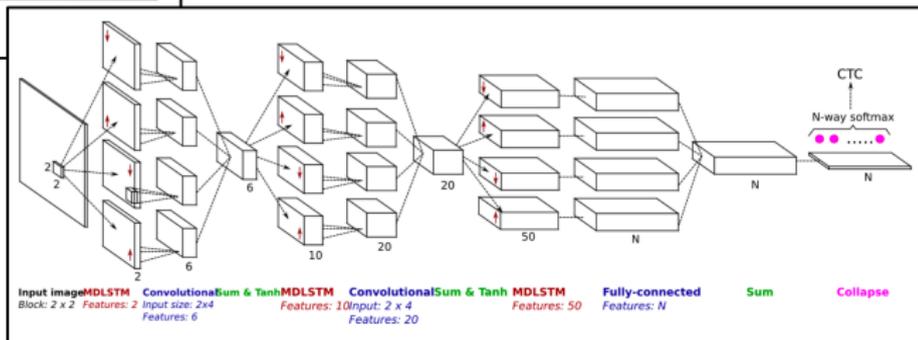
- Multi-Dimensional Long Short-Term Memory Recurrent Neural Networks.
- Connectionist Temporal Classification.



## Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks

**Alex Graves**  
TU Munich, Germany  
graves@in.tum.de

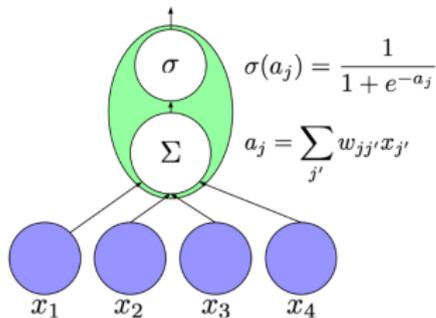
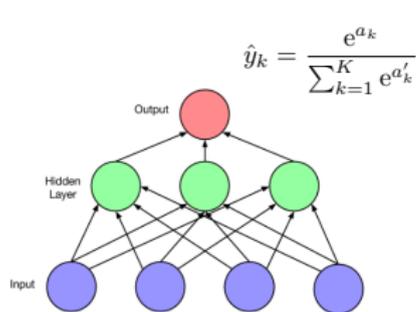
**Jürgen Schmidhuber**  
IDSIA, Switzerland and TU Munich, Germany  
juergen@idsia.ch



# Index

- 1 Motivation
- 2 Recurrent Neural Networks for Supervised Sequence Learning
- 3 Connectionist Temporal Classification
- 4 A DL network for Handwritten Text Recognition

# Neural Networks



$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K y_k^{(n)} \log(y_k^{(n)}) + (1 - y_k^{(n)}) \log(1 - y_k^{(n)})$$

Training: Back-propagation combined with Stochastic Gradient Descent

# The effective capacity of NNs is sufficient for memorizing the entire data set

## UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

**Chiyuan Zhang\***  
Massachusetts Institute of Technology  
chiyuan@mit.edu

**Samy Bengio**  
Google Brain  
bengio@google.com

**Moritz Hardt**  
Google Brain  
mrtz@google.com

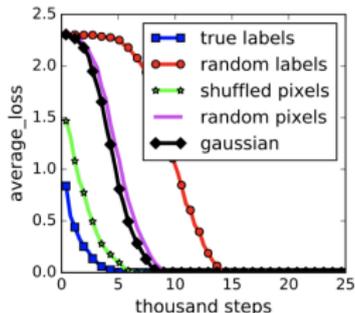
**Benjamin Recht†**  
University of California, Berkeley  
brecht@berkeley.edu

**Oriol Vinyals**  
Google DeepMind  
vinyals@google.com

Table 1: The training and test accuracy (in percentage) of various models on the CIFAR10 dataset. Performance with and without data augmentation and weight decay are compared. The results of fitting random labels are also included.

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75
(fitting random labels)		no	no	100.0	9.78
Inception w/o BatchNorm	1,649,402	no	yes	100.0	83.00
		no	no	100.0	82.00
		(fitting random labels)	no	no	100.0
Alexnet	1,387,786	yes	yes	99.90	81.22
		yes	no	99.82	79.66
		no	yes	100.0	77.36
		no	no	100.0	76.07
		(fitting random labels)	no	no	99.82
MLP 3x512	1,735,178	no	yes	100.0	53.35
		no	no	100.0	52.39
		(fitting random labels)	no	no	100.0
MLP 1x512	1,209,866	no	yes	99.80	50.39
		no	no	100.0	50.51
		(fitting random labels)	no	no	99.34

*“Deep NNs easily fit random labels”*



Fitting random labels and random pixels on CIFAR10

*“even with all of the regularizers turned off, all of the models still generalize very well”*

*“It is unlikely that the regularizers are the fundamental reason for generalization, as the networks continue to perform well after all the regularizers removed.”*

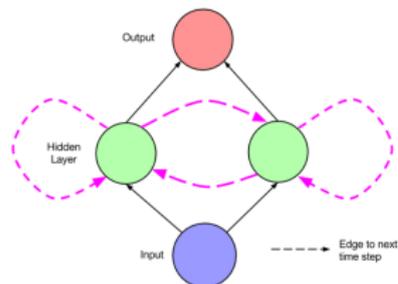
## Exploiting temporal correlation: Recurrent NNs

- Recurrent neural networks (RNNs) are connectionist models with the ability to selectively pass information across sequence steps, while processing sequential data one element at a time.
- Feedforward neural networks augmented by the inclusion of edges that span adjacent time steps, introducing a notion of time to the model.

$$\mathbf{h}^{(t)} = \sigma(W^{\text{hx}}\mathbf{x}^{(t)} + W^{\text{hh}}\mathbf{h}^{(t-1)} + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W^{\text{yh}}\mathbf{h}^{(t)} + \mathbf{b}_y)$$

$$\hat{\mathbf{y}}^{(t)} = P(\text{label}|\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$



Figures: Lipton, Berkowitz (2015) *A Critical Review of Recurrent Neural Networks for Sequence Learning*

## Why not Markov models?

- Generative models, EM-based training.
  - Markov model approaches are limited because their states must be drawn from a modestly sized discrete state space  $\mathcal{S}$
  - Standard operations become infeasible with an HMM when the set of possible hidden states grows large.
- 
- In contrast, RNNs can effectively handle a number of distinct states grows exponentially with the number of nodes in the layer.
  - Discriminative supervised models, SGD training.

## Bi-directional RNNs

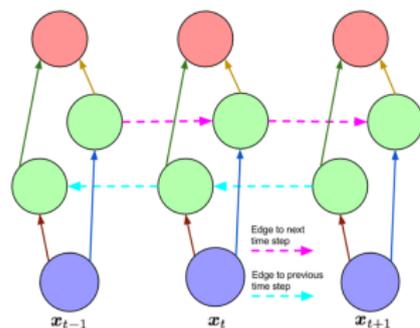
- There are two layers of hidden nodes. Both hidden layers are connected to input and output.
- The first layer has recurrent connections from the past time steps while in the second the direction of recurrent of connections is flipped.

$$\mathbf{h}^{(t)} = \sigma(W^{\text{hx}}\mathbf{x}^{(t)} + W^{\text{hh}}\mathbf{h}^{(t-1)} + \mathbf{b}_h)$$

$$\mathbf{z}^{(t)} = \sigma(W^{\text{zx}}\mathbf{x}^{(t)} + W^{\text{zz}}\mathbf{z}^{(t+1)} + \mathbf{b}_z)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(W^{\text{yh}}\mathbf{h}^{(t)} + W^{\text{yz}}\mathbf{z}^{(t)} + \mathbf{b}_y)$$

$$\hat{\mathbf{y}}^{(t)} = P(\text{label}|\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)})$$



Figures: Lipton, Berkowitz (2015) *A Critical Review of Recurrent Neural Networks for Sequence Learning*

## Long Short-Term Memory Cells

- RNNs are limited by their temporal expressiveness, which is manifested by the *vanishing gradient problem*. Don't have the ability to learn long-range dependencies.
- This problem is alleviated by LSTMs, which show superior ability to learn long-range dependencies.

$$\mathbf{g}^{(t)} = \phi(W^{\mathbf{g}\mathbf{x}}\mathbf{x}^{(t)} + W^{\mathbf{g}\mathbf{h}}\mathbf{h}^{(t-1)} + \mathbf{b}_g)$$

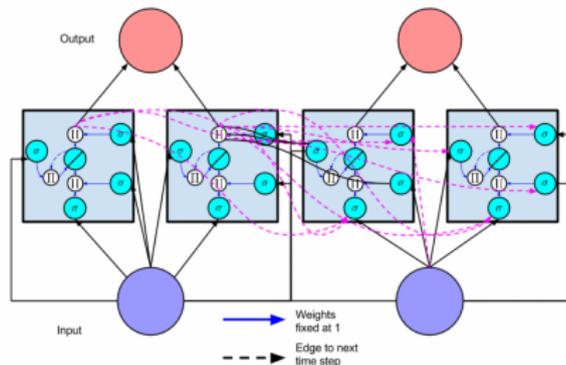
$$\mathbf{i}^{(t)} = \sigma(W^{\mathbf{i}\mathbf{x}}\mathbf{x}^{(t)} + W^{\mathbf{i}\mathbf{h}}\mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}^{(t)} = \sigma(W^{\mathbf{f}\mathbf{x}}\mathbf{x}^{(t)} + W^{\mathbf{f}\mathbf{h}}\mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{o}^{(t)} = \sigma(W^{\mathbf{o}\mathbf{x}}\mathbf{x}^{(t)} + W^{\mathbf{o}\mathbf{h}}\mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{s}^{(t)} = \mathbf{g}^{(t)} \odot \mathbf{i}^{(t)} + \mathbf{s}^{(t-1)} \odot \mathbf{f}^{(t)}$$

$$\mathbf{h}^{(t)} = \phi(\mathbf{s}^{(t)}) \odot \mathbf{o}^{(t)}$$



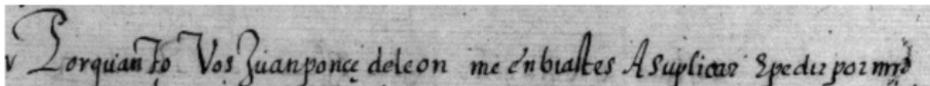
Figures: Lipton, Berkowitz (2015) *A Critical Review of Recurrent Neural Networks for Sequence Learning*

# Index

- 1 Motivation
- 2 Recurrent Neural Networks for Supervised Sequence Learning
- 3 Connectionist Temporal Classification**
- 4 A DL network for Handwritten Text Recognition

## Labelling alignment is a problem with RNNs

- Labelling unsegmented sequence data is a ubiquitous problem in real-world sequence learning.
- RNNs can only be trained to make a series of independent label classifications.
- Training data must be pre-segmented.
- Since the RNN network only outputs local classifications, a post-processing stage is required to give the final label sequence.



Porquanto Vos Juan ponce de leon me enbiastes A suplicar e pedir por merced

## Hybrid HMM-RNN models

- HMM models the sequential structure of the data, able to segment localized classifications provided by the RNNs into a temporal classification of the entire label sequence.
- In HMMs, training is generative, even though sequence labelling is a discriminative task.
- Iterative approach, where the alignment provided by the HMM is used to successively retrain the neural network.

Markovian Models for Sequential Data

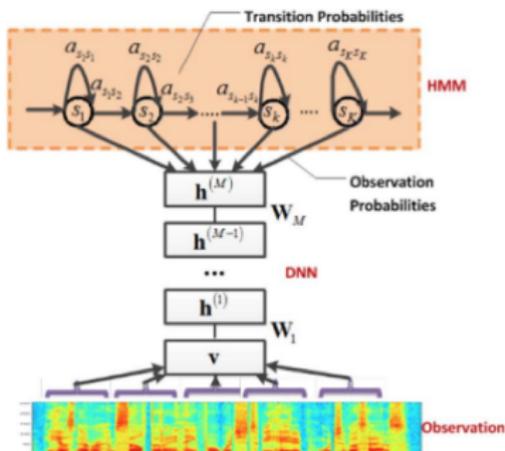
---

Yoshua Bengio †

Dept. Informatique et Recherche Opérationnelle  
Université de Montréal, Montreal, Qc H3C-3J7, Canada

# Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition

George E. Dahl, Dong Yu, Senior Member, IEEE, Li Deng, Fellow, IEEE, and Alex Acero, Fellow, IEEE

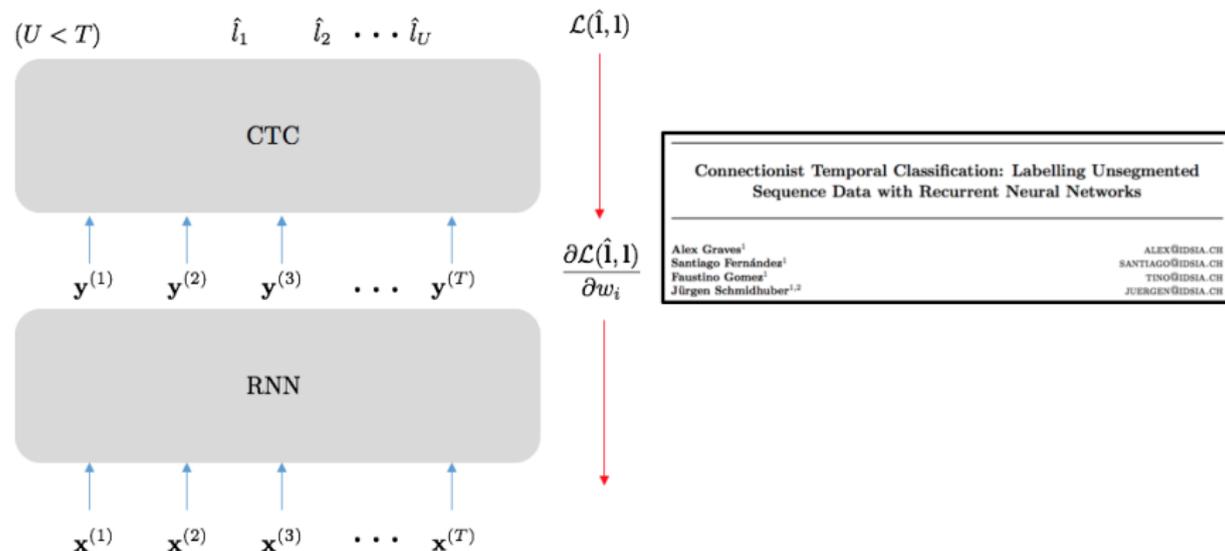


## Algorithmic 1 Main Steps to Train CD-DNN-HMMs

- 1) Train a best tied-state CD-GMM-HMM system where state tying is determined based on the data-driven decision tree. Denote the CD-GMM-HMM  $gmm-hmm$ .
- 2) Parse  $gmm-hmm$  and give each senone name an ordered *senoneid* starting from 0. The *senoneid* will be served as the training label for DNN fine-tuning.
- 3) Parse  $gmm-hmm$  and generate a mapping from each physical tri-phone state (e.g., b-ah+Ls2) to the corresponding *senoneid*. Denote this mapping *state2id*.
- 4) Convert  $gmm-hmm$  to the corresponding CD-DNN-HMM  $dnn-hmm1$  by borrowing the tri-phone and senone structure as well as the transition probabilities from  $gmm-hmm$ .
- 5) Pre-train each layer in the DNN bottom-up layer by layer and call the result *ptdnn*.
- 6) Use  $gmm-hmm$  to generate a state-level alignment on the training set. Denote the alignment *align-raw*.
- 7) Convert *align-raw* to *align* where each physical tri-phone state is converted to *senoneid*.
- 8) Use the *senoneid* associated with each frame in *align* to fine-tune the DBN using back-propagation or other approaches, starting from *ptdnn*. Denote the DBN *dnn*.
- 9) Estimate the prior probability  $p(s_i) = n(s_i)/n$ , where  $n(s_i)$  is the number of frames associated with senone  $s_i$  in *align* and  $n$  is the total number of frames.
- 10) Re-estimate the transition probabilities using *dnn* and  $dnn-hmm1$  to maximize the likelihood of observing the features. Denote the new CD-DNN-HMM  $dnn-hmm2$ .
- 11) Exit if no recognition accuracy improvement is observed in the development set; Otherwise use

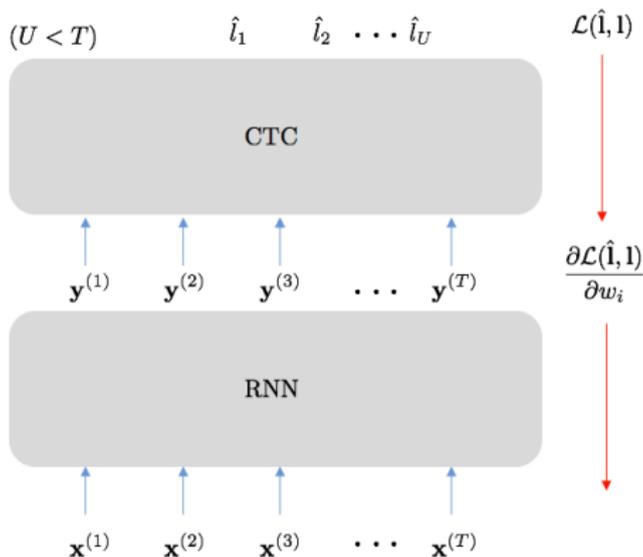
## Connectionist Temporal Classification (CTC)

- Pure discriminative methods tend to give better results for classification tasks, because they focus entirely on finding the correct labels.
- CTC models all aspects of the sequence with a single neural network.
- It also does not require pre-segmented training data, or external post-processing to extract the label sequence from the network outputs.



## Notation

- $\mathcal{A}$  is the label alphabet. E.g., characters in the Latin alphabet.
- $\mathcal{A}' = \mathcal{A} \cup \{\text{blank}\}$
- $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(T)})$  is the input sequence
- $\mathbf{Y} = (\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(T)})$  is the RNN output, where  $y_k^{(t)} \in [0, 1]$ ,  $k = 1, \dots, |\mathcal{A}'|$ , and  $\sum_{k=1}^{|\mathcal{A}'|} y_k^{(t)} = 1$  (softmax layer at the RNN output).
- $\hat{\mathbf{l}} = (\hat{l}_1, \hat{l}_2, \dots, \hat{l}_U)$  is the out label sequence,  $\hat{l}_j \in \mathcal{A}'$ .



## From Network Outputs to Labellings

- $y_k^{(t)}$  → Probability of observing label  $k$  at time  $t$ .
- $\mathbf{Y}$  → Probability distribution over all possible *paths* in the set  $\tilde{\mathcal{A}}^T$ :

$$p(\boldsymbol{\pi}|\mathbf{X}) = \prod_{i=1}^T y_{\pi_t}^{(t)} \quad \boldsymbol{\pi} \in \tilde{\mathcal{A}}^T$$

- A many-to-one map  $\mathcal{F} : \tilde{\mathcal{A}}^T \rightarrow \tilde{\mathcal{A}}^{\leq T}$  is defined.
  - ▶ Remove repeated labels
  - ▶ Remove blanks
  - ▶  $\mathcal{F}(a - ab-) = \mathcal{F}(-aa - -abb) = aab$
  - ▶ Outputting a new label when the network switches from predicting no label to predicting a label, or from predicting one label to another.
- For any  $I \in \mathcal{A}^U$  with  $U \leq T$  then

$$p(I|\mathbf{X}) = \sum_{\boldsymbol{\pi} \in \mathcal{F}^{-1}(I)} p(\boldsymbol{\pi}|\mathbf{X})$$

## From Network Outputs to Labellings

- $y_k^{(t)}$  → Probability of observing label  $k$  at time  $t$ .
- $\mathbf{Y}$  → Probability distribution over all possible *paths* in the set  $\tilde{\mathcal{A}}^T$ :

$$p(\boldsymbol{\pi}|\mathbf{X}) = \prod_{i=1}^T y_{\pi_t}^{(t)} \quad \boldsymbol{\pi} \in \tilde{\mathcal{A}}^T$$

- A many-to-one map  $\mathcal{F} : \tilde{\mathcal{A}}^T \rightarrow \tilde{\mathcal{A}}^{\leq T}$  is defined.
  - ▶ Remove repeated labels
  - ▶ Remove blanks
  - ▶  $\mathcal{F}(a - ab-) = \mathcal{F}(-aa - -abb) = aab$
  - ▶ Outputting a new label when the network switches from predicting no label to predicting a label, or from predicting one label to another.
- For any  $l \in \mathcal{A}^U$  with  $U \leq T$  then

$$p(l|\mathbf{X}) = \sum_{\boldsymbol{\pi} \in \mathcal{F}^{-1}(l)} p(\boldsymbol{\pi}|\mathbf{X})$$

## From Network Outputs to Labellings

- $y_k^{(t)}$  → Probability of observing label  $k$  at time  $t$ .
- $\mathbf{Y}$  → Probability distribution over all possible *paths* in the set  $\tilde{\mathcal{A}}^T$ :

$$p(\boldsymbol{\pi}|\mathbf{X}) = \prod_{i=1}^T y_{\pi_t}^{(t)} \quad \boldsymbol{\pi} \in \tilde{\mathcal{A}}^T$$

- A many-to-one map  $\mathcal{F} : \tilde{\mathcal{A}}^T \rightarrow \tilde{\mathcal{A}}^{\leq T}$  is defined.
  - ▶ Remove repeated labels
  - ▶ Remove blanks
  - ▶  $\mathcal{F}(a - ab-) = \mathcal{F}(-aa - -abb) = aab$
  - ▶ Outputting a new label when the network switches from predicting no label to predicting a label, or from predicting one label to another.
- For any  $I \in \mathcal{A}^U$  with  $U \leq T$  then

$$p(I|\mathbf{X}) = \sum_{\boldsymbol{\pi} \in \mathcal{F}^{-1}(I)} p(\boldsymbol{\pi}|\mathbf{X})$$

## From Network Outputs to Labellings

$$p(I|\mathbf{X}) = \sum_{\pi \in \mathcal{F}^{-1}(I)} p(\pi|\mathbf{X}), \quad I \in \mathcal{A}^U$$

- This collapse of different paths into the same label sequence  $I$  is what makes possible for CTC to use unsegmented data.
- $p(I|\mathbf{X})$  can be efficiently computed using dynamic programming  $\rightarrow$  **CTC Forward-Backward (CTC-FB) algorithm**.
- Further, the gradient of  $-\log p(I|\mathbf{X})$  w.r.t.  $y_k^t$  can be computed using **intermediate messages** computed during CTC-FB.

## CTC Forward-Backward (I)

- Given  $I$ , we define  $I'$  with blanks added to the beginning and end and inserted every pair of labels.
- For  $t = 1, \dots, T$  and  $u = 1 \dots, 2U + 1$  we define the following two sets of paths and cumulated probabilities

$$V(t, u) = \left\{ \pi \in \tilde{\mathcal{A}}^t : \mathcal{F}(\pi) = I_{\lfloor \frac{u}{2} \rfloor}, \pi_t = I'_u \right\} \rightarrow \alpha(t, u) = \sum_{\pi \in V(t, u)} \prod_{i=1}^t y_{\pi_i^{(i)}}$$

$$W(t, u) = \left\{ \pi \in \tilde{\mathcal{A}}^{T-t} : \mathcal{F}(\hat{\pi} + \pi) = I \quad \forall \hat{\pi} \in V(t, u) \right\} \rightarrow \beta(t, u) = \sum_{\pi \in W(t, u)} \prod_{i=t}^T y_{\pi_i^{(i)}}$$

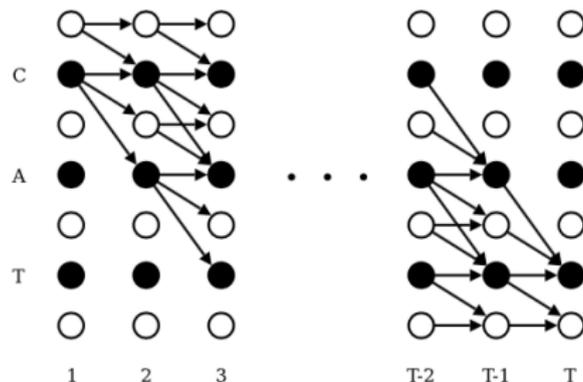


Figure: A.Graves et al. *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. (2006)

## CTC Forward-Backward (II)

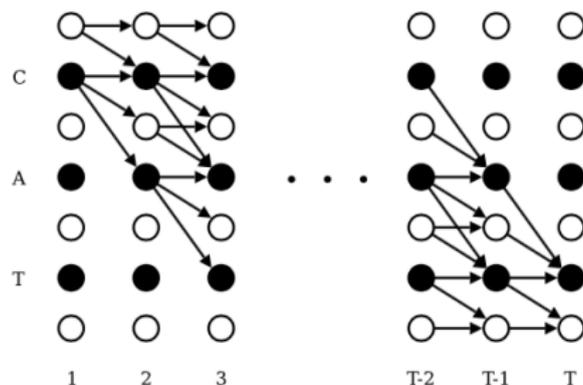


Figure: A.Graves et al. *Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks*. (2006)

$$\alpha(t, u) = f\left(\alpha(t-1, u), \alpha(t-1, u-1), \alpha(t-1, u-2)\right)$$

$$\beta(t, u) = g\left(\beta(t+1, u), \beta(t+1, u+1), \alpha(t+1, u+2)\right)$$

$$p(I|\mathbf{X}) = \alpha(t, 2U) + \alpha(t, 2U + 1)$$

$$\frac{-\log p(I|\mathbf{X})}{\partial y_k^{(t)}} = -\frac{1}{p(I|\mathbf{X})y_k^{(t)}} \sum_{u: I'_u=k} \alpha(t, u)\beta(t, u)$$

## CTC decoding

Upon training, for a new input sequence  $\mathbf{X}$ , we would ideally compute

$$I^* = \arg \max_I p(I|\mathbf{x})$$

### Best Path Decoding (Trivial)

Let  $\pi^* = \arg \max_{\pi} p(\pi|\mathbf{X})$ , then

$$I^* \approx \mathcal{F}(\pi^*)$$

### Prefix Search Decoding

Modified Forward/Backward algorithm to calculate the probabilities of successive extensions to labelling prefixes. However, the number of prefixes it must expand grows exponentially with the input sequence length.

### Constrained Decoding

Constrain the output labellings according to some predefined grammar

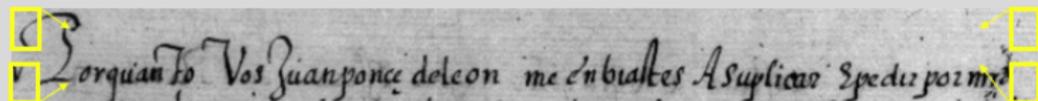
# Index

- 1 Motivation
- 2 Recurrent Neural Networks for Supervised Sequence Learning
- 3 Connectionist Temporal Classification
- 4 A DL network for Handwritten Text Recognition

## Multi-dimensional LSTMs

- Generalize bi-RNNs by providing recurrent connections along all spatio-temporal dimensions present in the data.
- Robust to local distortions along any combination of input dimensions.

Porquanto Vos Juan ponce de leon me enbiastes A suplicar e pedir por merced



$$h_{i,j}^{NW} = \text{LSTM} \left( x_{i,j}, h_{i,j-1}^{NW}, h_{i-1,j}^{NW}, h_{i-1,j+1}^{NW} \right)$$

$$h_{i,j}^{NE} = \text{LSTM} \left( x_{i,j}, h_{i,j+1}^{NE}, h_{i-1,j}^{NE}, h_{i-1,j+1}^{NE} \right)$$

$$h_{i,j}^{SW} = \text{LSTM} \left( x_{i,j}, h_{i,j+1}^{SW}, h_{i+1,j}^{SW}, h_{i+1,j+1}^{SW} \right)$$

$$h_{i,j}^{SE} = \text{LSTM} \left( x_{i,j}, h_{i,j+1}^{SE}, h_{i+1,j}^{SE}, h_{i+1,j+1}^{SE} \right)$$

## Network Hierarchy

- Subsampling layers, alternating between MLTSMs and feedforward/convolutional layers.
- The number of features computed by these layers increases as the size of the feature maps decrease.
- At the top of this network, there is one feature map for each label.
- A collapsing layer sums the features over the vertical axis, yielding a sequence of prediction vectors, effectively delaying the 2D to 1D transformation just before the character predictions, normalized with a softmax activation.

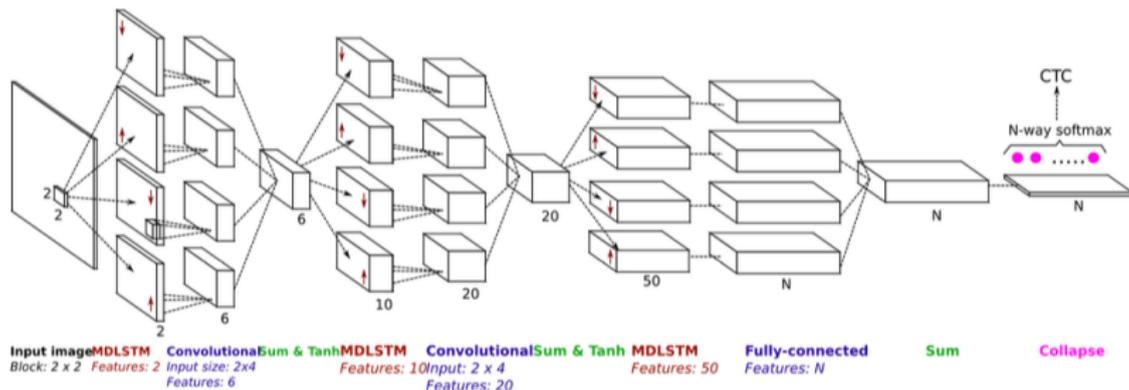


Figure: Pham et al. *Dropout improves recurrent neural networks for handwriting recognition* (2014)

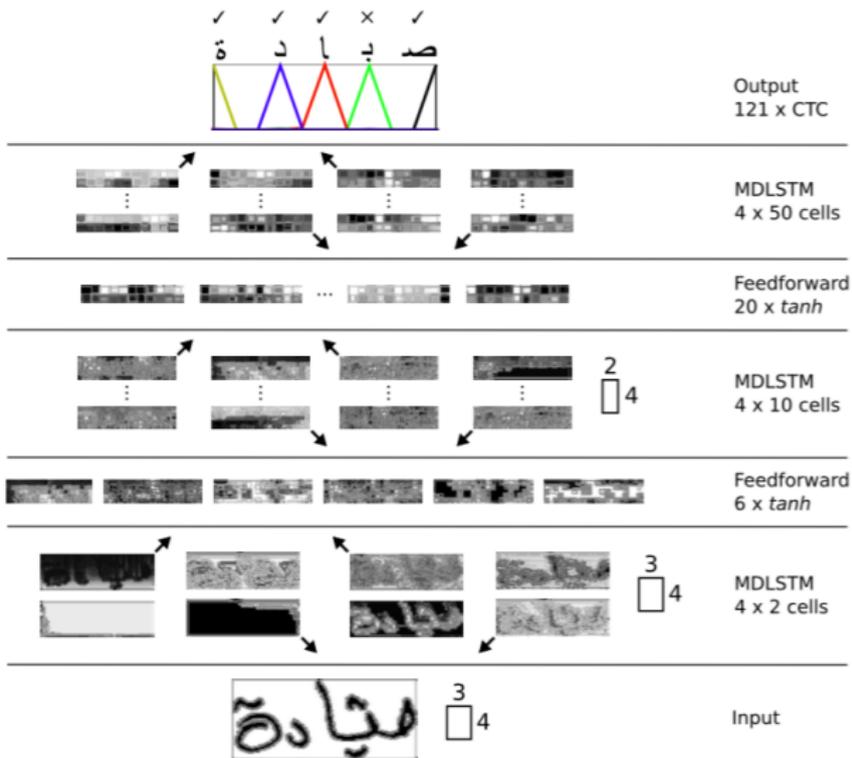


Figure: A.Graves and J.Schmidhuber. *Offline Handwriting Recognition with Multidimensional Recurrent Neural Networks.* (2008)