	rences
Importance of recurrent layers for unconstrained	
Denis Coquenet, Yann Soullard, Clément Chatelain, Thierry	J

LITIS Laboratory - EA 4108 Normandie University - University of Rouen, France

SIFED, 6<sup>th</sup> June 2019













#### Constraints

- Images (input) of variable size
- Sequence of characters (output) of variable length



• Towards a heavy use of neural networks



Connectionist Temporal Classification (CTC)

• Focus on optical model only without language model nor lexicon constraints



• Recurrence over horizontal and vertical axis (in both directions) : 4 LSTM/layer



• Recurrence over horizontal axis only (in both directions) : 2 LSTM/layer

#### Non-recurrent models - Convolutionnal Neural Network (CNN)



## State of the art

#### CNN for handwriting recognition

- Fully Convolutional Networks (FCN) + CTC [Ptucha2018]
  - Standard CNN without dense layer
- FCN with gating mechanism + CTC [Yousef2018]
  - Gates (tanh, sigmoid)
  - Residual connections
  - Depthwise Separable Convolutions
  - High normalization (batch & layer)
- FCN with gating mechanism + CTC[Ingle2019]
  - Gates (ReLU, sigmoid)
  - Shared weight layers

Context	Studied architectures	Experiments	Conclusion	References
Context				

#### Objective

Design a convolutional network competitive with recurrent ones to reduce the training time :  $\ensuremath{\mathsf{G-CNN}}$ 

#### Main questions

- Are recurrent layers really necessary for handwriting recognition ?
- Are CNN really lighter than recurrent model in terms of parameters ?
- Is it possible to easily obtain competitive results without recurrence ?



## Our baseline model - CNN+BLSTM



#### Features

- From [Soullard2019] (state-of-the-art results)
- Recurrent model
- 8 convolutions
- 2.5 million of parameters

n : number of characters in the alphabet



## Our G-CNN - gates

#### Gating mechanism



## Training details

#### Hyperparameters

- Sliding window : 32x32 px
- Loss : CTC
- Optimizer : Adam
- Initial learning rate :  $10^{-4}$
- Momentum : 0.9

#### Criteria

- Character Error Rate (CER)
- Number of parameters
- Training time

#### Raw model comparison

- We focus only on the network performance alone
- No language model
- No lexicon constraints

#### Dataset

RIMES (lines)

Conte	ext Studied architectures	Experiments	Conclusion	References
RII	MES dataset			
	Dataset characteristics			
	<ul> <li>+1,300 writers</li> </ul>			

- French writings
- 12,723 pages segmented into lines

## RIMES dataset split

Training	Validation	Test	Alphabet
9,947	1,333	778	100

## Example

## First experiment : Raw comparison

Architecture	CER(%)	CER (%)	Training	Parameters (M)
Architecture	validation	test	time	r arameters (m)
CNN+BLSTM	6.98	6.88	1d22h59	4.1
CNN+Dense only	17.73	19.03	1h10	1.5
G-CNN	9.92	10.03	10h00	6.9

• BLSTM layers responsible for a large amount of parameters (2.6 M)

## First experiment : Raw comparison

Architecture	CER(%)	CER (%)	Training	Parameters (M)
Architecture	validation	test	time	r arameters (m)
CNN+BLSTM	6.98	6.88	1d22h59	4.1
CNN+Dense only	17.73	19.03	1h10	1.5
G-CNN	9.92	10.03	10h00	6.9

- BLSTM layers responsible for a large amount of parameters (2.6 M)
- BLSTM layers increase performance dramatically (-12.15% in test)

## First experiment : Raw comparison

Architecture	CER(%)	CER (%)	Training	Parameters (M)
Architecture	validation	test	time	r arameters (m)
CNN+BLSTM	6.98	6.88	1d22h59	4.1
CNN+Dense only	17.73	19.03	1h10	1.5
G-CNN	9.92	10.03	10h00	6.9

- BLSTM layers responsible for a large amount of parameters (2.6 M)
- BLSTM layers increase performance dramatically (-12.15% in test)
- G-CNN : more parameters but training time shorter (parallel computing)

## Second experiment - Robustness against complexified data

#### Modified version of RIMES dataset

Lined paper background addition



## Second experiment - Robustness against complexified data

Architactura	Rackground	CER(%)	CER (%)	Training
Architecture	Dackground	validation	test	time
	Without	6.98	6.88	1d22h59
	With	8.81	9.27	1d1h29
	Without	9.92	10.03	10h00
G-CININ	With	11.70	12.55	8h27

• Similar behavior - CER increased by 2.39% for the CNN+BLSTM and 2.52% for the G-CNN (in test)

## Third experiment - Impact of data augmentation

	RIMES
1. Raw	Par la præzonte je vaue fais part de ma
2. Contrast	Par la præsente je vaus fais part de ma
3. Sign flipping	Par la præsate je vaue faie part de ma
4. Long scaling	Par la présente je vaue faie part de ma
5. Short scaling	Par la presente je vaue faie part de ma
6. Width dilation	Par la præsente je vaue fale part de ma
7. Height dilation	Par la présente je voue fais part de ma

Importance of recurrent layers for unconstrained handwriting recognition

Architecture	Data augmentation	CER(%) - validation	CER (%) - test
	Without	6.98	6.88
CININ+DL3 I M	With	6.59	5.94
	Without	9.92	10.03
G-CININ	With	8.93	8.73

- CER decreased by 1.30% for the G-CNN and 0.94% for the CNN+BLSTM
- Assumption : G-CNN needs more examples whereas CNN+BLSTM compensates with its use of context

Architactura	CER(%)	CER (%)	Training	Daramatars (M)
Arciitecture	validation	test	time	Farameters (W)
G-CNN	9.92	10.03	10h00	6.9
(1) Only standard convolutions	10.02	9.97	6h41	9.0
(2) Max pooling from the very beginning	13.31	13.35	2h57	6.9
(3) No shared weight layers	9.78	9.85	8h54	7.7

• (1) Depthwise Separable Convolutions enables saving 2.1 M of parameters preserving the performance (+0.06%)

Architecture	CER(%) validation	CER (%) test	Training time	Parameters (M)
G-CNN	9.92	10.03	10h00	6.9
(1) Only standard convolutions	10.02	9.97	6h41	9.0
(2) Max pooling from the very beginning	13.31	13.35	2h57	6.9
(3) No shared weight layers	9.78	9.85	8h54	7.7

- (1) Depthwise Separable Convolutions enables saving 2.1 M of parameters preserving the performance (+0.06%)
- (2) Delaying the use of max pooling increases the performance (by 3.32%)

Architactura	CER(%)	CER (%)	Training	Daramatars (M)	
	validation	test	time	Farameters (W)	
G-CNN	9.92	10.03	10h00	6.9	
(1) Only standard convolutions	10.02	9.97	6h41	9.0	
(2) Max pooling from the very beginning	13.31	13.35	2h57	6.9	
(3) No shared weight layers	9.78	9.85	8h54	7.7	

- (1) Depthwise Separable Convolutions enables saving 2.1 M of parameters preserving the performance (+0.06%)
- (2) Delaying the use of max pooling increases the performance (by 3.32%)
- (3) Shared weight layers enable saving 0.8 M parameters, with a similar CER (+0.18% only)

Architecture	CER(%)	CER (%)	Training	Parameters (M)	
Alcintecture	validation	test	time		
G-CNN	9.92	10.03	10h00	6.9	
(4) Doubled convolutions in GateBlocks	9.96	10.15	4h10	7.4	
(5) Removal of the 2 GateBlocks	10.09	10.33	6h37	6.1	

• (4) Increasing the number of convolutions between gates is not necessary (+0.12%)

Architecture	CER(%) validation	CER (%) test	Training time	Parameters (M)
G-CNN	9.92	10.03	10h00	6.9
(4) Doubled convolutions in GateBlocks	9.96	10.15	4h10	7.4
(5) Removal of the 2 GateBlocks	10.09	10.33	6h37	6.1

- (4) Increasing the number of convolutions between gates is not necessary (+0.12%)
- (5) The majority of the work is done before the GateBlocks (+0.3%)

## Conclusion - recurrent models

#### Structure

- Convolutional part (CNN): feature extraction
- Recurrent part (LSTM): sequence modeling

#### Advantages

- Performance
- Simple architectures

#### Drawbacks

Recurrent models have long training times:

- LSTM implies a large amount of parameters
- Recurrence implies sequential computations

# Conclusion - G-CNN models

#### Structure

- Feature extraction similar to CNN+BLSTM
- Gating mechanism to filter information

#### Advantages

Convolution = parallelizable operation & few parameters

- Reduced training time
- Deeper networks, bigger receptive fields

#### Drawbacks

- Number of hyperparameters, hard tuning
- Complex architecture
- Performance hardly competitive

#### Toward an even lighter network

# Give up densely connected layers to build a Fully Convolutionnal Network

#### Exploring other alternatives

- Attention models [Chowdhury2018]
- Dense net [Huang2016]



Context	Studied architectures	Experiments	Conclusion	References
References	;			
[Pham2014	] V. Pham et Handwriting	al. "Dropout Improves R Recognition". In: <i>ICFHI</i>	ecurrent Neural Net R (2014).	works for
[Huang201	6] Huang et al.	Densely Connected Con	volutional Networks	. 2016.
[Puigcerver	2017] J. Puigcerve Necessary fo pp. 67–72.	r. "Are Multidimensional r Handwritten Text Reco	Recurrent Layers Rognition?". In: ICDA	eally NR. 2017,
[Chowdhur	y2018] C. Arindam Handwritten	et al. An Efficient End-to Text Recognition. 2018	o-End Neural Model	for
[Yousef201	8] M. Yousef e Recognition	t al. Accurate, Data-Effic with Convolutional Neur	cient, Unconstrained al Networks. 2018.	Text
[Ptucha201	.8] Felipe Petro using Fully ( <i>Recognition</i>	ski Such et al. "Intelligen Convolutional Neural Net 88 (Dec. 2018).	t Character Recogni works". In: <i>Pattern</i>	ition
[Soullard20	19] Y. Soullard Temporal Cl	et al. CTCModel: a Kera assification. 2019.	s Model for Connect	tionist
[Ingle2019]	R. Ingle et a 2019.	I. A Scalable Handwritte	n Text Recognition	System.