

# Textual Data Selection for Language Modelling in the Scope of Automatic Speech Recognition

F. Mezzoudj<sup>1</sup>, D. Langlois<sup>2</sup>, D. Jouvet<sup>2</sup>, A. Benyettou<sup>1</sup>

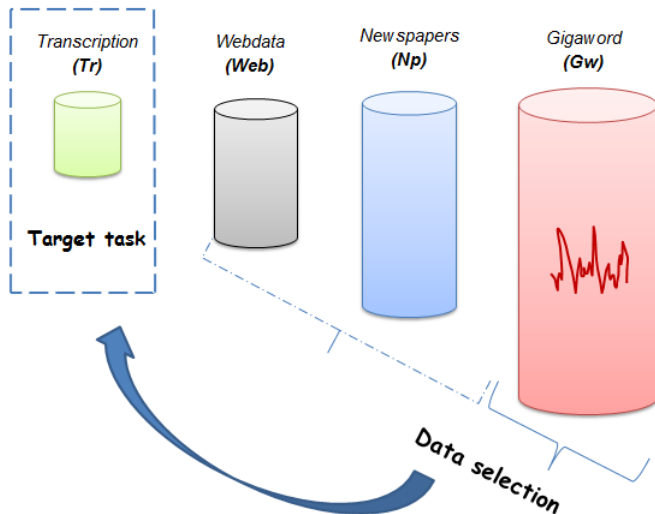
<sup>1</sup> USTO/UHBC-Algeria, <sup>2</sup> INRIA, LORIA, France

ICNLSP, Algiers, October 18<sup>th</sup>, 2015

# Introduction

- **Language model** (LM) is an important module in ASR systems.
- Learning LM requires a large amount of textual data.
- A high-performance LM is trained using a small corpus close to the target task (**in-domain**) and a huge corpus not close to this task (**non-domain**).
- We investigate selection of French textual data in order to improve LMs for automatic speech transcription of **broadcast news** & **TV shows**.

# Introduction



# Outline

- 1 Introduction
- 2 Data Selection
- 3 Experimental Setup
- 4 Data selection strategy
- 5 Transcription experiments
- 6 Conclusion

# Data selection / Related work

## Application on 2 sources :

- **Klakow (2000)** used a **log-likelihood** criterion to select newspaper articles.
- **Wang et al.(2002)** selected text units from the non-domain with lowest **perplexity** according to the in-domain LM.
- **Moore et al. (2010)** selected sentences from the non-domain with lowest **difference cross- entropy** according to 2 LMs with equal size, representing the in-domain LM and a non-domain LM).

# Data selection / New situation

We face a completely different situation :

- We use **4** corpora corresponding to different sources :
  - **manual transcriptions** of broadcast news & TV shows ;
  - **Webdata** (from Web sites : Magazines, TV) ;
  - **Newspapers** (Le Monde & L'Humanité) ;
  - **Gigaword** corpus 2<sup>nd</sup> edition.
- Each corpus contributes differently to the final LM's training.
- Corpora may be noisy because of the variable quality of the sources.

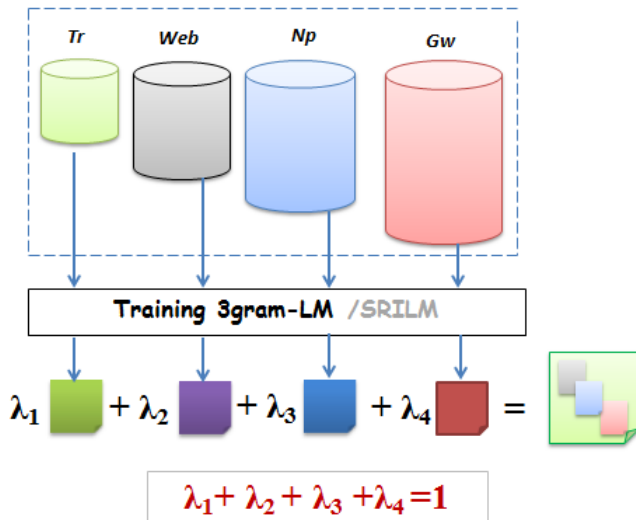
# Experimental setup

- **Training corpora** Sources & # words [M]

$Tr$ (radio broadcast transcriptions)	114
$Web$ (web data)	334
$Np$ (newspapers)	526
$Gw$ (gigaword corpus)	783
$Tr + Web + Np + Gw$	<b>1 757</b>

- **Validation Corpus**,  $DevLM \simeq 300$  K words
- **Test Corpus**,  $TestLM \simeq 90$  K words
- **Vocabulary**  $\simeq 100$ K words

# Experimental setup / baseline LM





# Experimental setup / baseline LM

Table : **Baseline LM**, interpolated from the individual LMs.

Sources	Interpolated LM		
	weights	ppl <i>DevLM</i>	ppl <i>TestLM</i>
<i>Tr</i>	0.685	185.7	<b>218.9</b>
<i>Web</i>	0.246		
<i>Np</i>	0.062		
<b>GW</b>	0.007		

- large difference in the weight of the individual LMs ;
- Gw brings small contribution in the final LM.

# Data selection strategy

## Step 1

- # source used :

2

- data selection

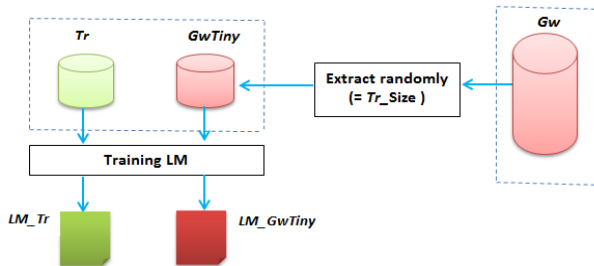
on : **Gw**

- "in-domain" :

*LM\_Tr*

- "non-domain" :

*LM\_GwTiny*

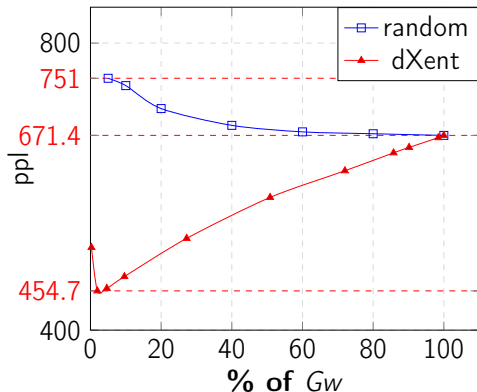


## Step 2



# Data selection strategy

- The LM's ppl obtained by using the whole Gw corpus is 671.4;
- Small subsets selected **randomly** on the Gw data **degrades** the ppl;
- The difference cross-entropy (**dXent**) selection data on the Gw corpus **improves** the ppl.



# Multisource-LM / approach 1

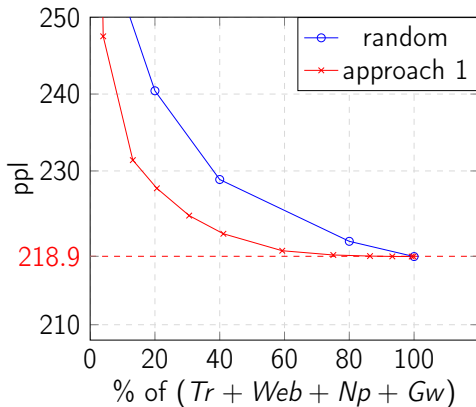
## approach 1

- # source used : 4 ;
- selected data :  
*(Tr+ Web+ Np+ GW)*;
- "in-domain" : *LM\_Tr*;
- "non-domain" : *LM\_GwTiny*;



# Multisource-LM / approach 1

- The LM's ppl obtained by using the 4 corpora is 218.9
- The selection applied on data with a **random** process **degrades** the ppl;
- The selection applied on data with the ***dXent*** computed with  $LM_{(Tr)}$  and  $LM_{GwTiny}$  (approach 1) **doesn't improve** the ppl.



# Multisource-LM / approach 2

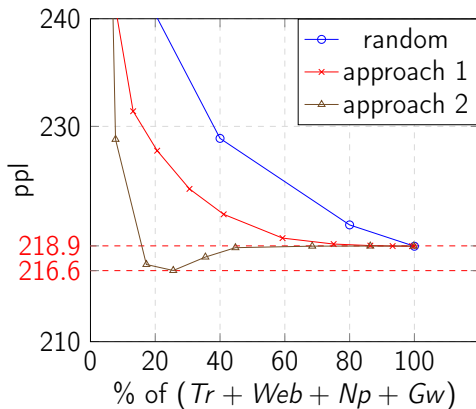
## approach 2

- # source used : 4 ;
- selected data :  
(*Tr* + *Web* + *Np* + *Gw*) ;
- "in-domain" : *LM\_TrWebNp* ;
- "non-domain" : *LM\_Gw* .



# Multisource-LM / approach 2

- The selection (approach 2) applied on  $(Tr, Web, Np, Gw)$  data with the ***dXent*** [computed with  $LM_{(TrWebNp)}$  and  $LM_{Gw}$ ] improves the ppl.





# Multisource-LM / approach 2

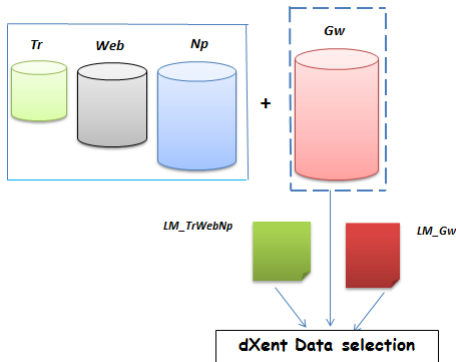
**Table :** Best LM, interpolated from the individual source LMs, after data selection using **approach 2**.

Sources	Interpolated LM		
	weights	ppl <i>DevLM</i>	ppl <i>TestLM</i>
<i>Tr</i> (88%)	0.608	185.1	<b>216.6</b>
<i>Web</i> (62%)	0.234		
<i>Np</i> (26%)	0.062		
<i>Gw</i> (0.2%)	<b>0.096</b>		

# Multisource-LM / approach 3

## approach 3

- # source used : 4 ;
- selected data : **GW** ;
- "in-domain" :  
*LM\_TrWebNp* ;
- "non-domain" : *LM\_Gw*.



# vs approach 2

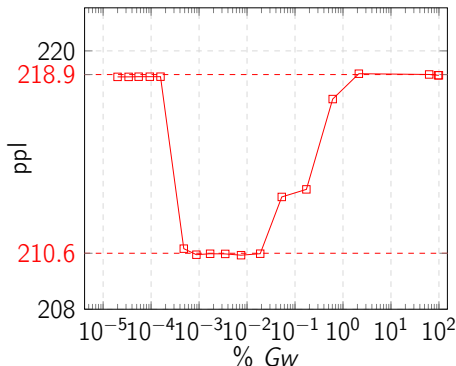
## approach 2

- # source used : 4 ;
- selected data :  
(*Tr* + *Web* + *Np* + *Gw*) ;
- "in-domain" : *LM\_TrWebNp* ;
- "non-domain" : *LM\_Gw*.



# Multisource-LM / approach 3

- The selection applied on  $G_w$  data [with the ***dXent*** scored by  $LM_{(TrWebNp)}$  and  $LM_{Gw}$ ] added to  $Tr$ ,  $Web$  and  $Np$  data (approach 3) **improves** the ppl.



# Multisource-LM / approach 3

**Table :** Best LM, interpolated from the individual source LMs, after data selection using **approach 3**.

Sources	Interpolated LM		
	weights	ppl <i>DevLM</i>	ppl <i>TestLM</i>
<i>Tr</i> (100%)	0.660	179.9	<b>210.6</b>
<i>Web</i> (100%)	0.240		
<i>Np</i> (100%)	0.065		
<i>Gw</i> (0.05%)	<b>0.054</b>		

# Transcription experiments

- The speech corpora used come from **ESTER2**, **ETAPE** evaluation campaigns and the **EPAC** project.
- The speech transcription system relies on a diarization step and on the **Sphinx toolkit**.
- **39 HTK MFCC** features are used (+ 1<sup>st</sup> & 2<sup>nd</sup> temporal derivatives).

# Transcription experiments

LM	Size (gz file)	Etape Dev corpus	
		ppl	WER[%]
$(Tr + Web + Np + Gw)$	1.2 Gb	218.9	27.84
$(Tr + Web + Np)$	809.8 Mb	218.9	27.82
LM(app.2, thresh. -0.3)	391.3 Mb	217.2	28.07
LM(app.2, thresh. -0.2)	501.6 Mb	216.6	27.89
LM(app.3, thresh. -0.6)	809.3 Mb	210.6	27.68

- The best LM is trained with 55.4% of (Tr,Web,Np,Gw).

# Conclusion

- The choice of the LMs that represent **in-domain** and **non-domain** is important for data selection.
- Keeping the 3 data sources ( $Tr$ ,  $Web$ ,  $Np$ ) and **selecting** data from the **Gw** corpus with the  $dXent$  leads to better results than when selecting data from the whole corpora ( $Tr$ ,  $Web$ ,  $Np$ ,  $Gw$ ).
- We obtain competitive results in WER with **reducing** strongly the **size** of training corpus for LMs.



# Perspective

- We have to explore other ways to select data from the huge Gw in order to improve the LM performance.
- The vocabulary is a crucial module for transcription, so we take into account the time period of data sub-parts.

# For Further Reading I



R.C. Moore and W. Lewis.

"Intelligent selection of language model training data".  
In Proceedings of the ACL 2010 Conference Short  
Papers. pp. 220-224



G. Gravier, et al.

*"The etape corpus for evaluation of speech-based tv  
content processing in the french language", in  
LREC-Eighth ICLRE, 2012, p. na.*

Thank you  
for your attention !