

# SUBWORD LATENT SEMANTIC ANALYSIS FOR TEXTTILING-BASED AUTOMATIC STORY SEGMENTATION OF CHINESE BROADCAST NEWS

Yulian Yang, Lei Xie

A, S L P G, S f C S  
 N v P U, X  
 @ v . . . , @ v . - .

## ABSTRACT

(LSA) 8 f ffi T -  
 / - f ff f f I  
 (ASR) , v LSA, f  
 f C v T T f  
 f v LSA  
 f W v LSA  
 f S v T T  
 C - f- (OOV) v  
 f LSA T v E  
 LSA- T T TDT2 M v v BN) f f f  
 f C - LSA- T T (LVCSR). S  
 F1- f 0.6598 v T T 6.5% f 17.4%  
 T T 30% f E -

**Index Terms—**

## 1. INTRODUCTION

A I f f BN  
 S fi f v (BN) M f BN  
 / 11, v C  
 f v f C BN 7. S v  
 T f v / v f C E -  
 1, 2 / v f C f C -  
 3, f M 4 f 5. L 7. I v LSA, f v W  
 v v T T 6 v v ff W  
 BN 7 f v f  
 v LSA

## 2. CORPUS

T v v N N S F -  
 f C (60802085), R F f D P f  
 H E C (20070699015), N S B R -  
 P f S P (2007F15), NPU A S P W  
 (07XE0150), NPU F f F R (W018103). 53 f VOA M C BN. T 177

TDT2 LVCSR  
 37%, 20%  
 W P  
 90 (1390)  
 87  
 A  
 15-  
 f

**3. LATENT SEMANTIC ANALYSIS**

(LSA)  
 12.H  
 PCA 10  
 LSA  
 $\Gamma = \{t_1, t_2, t_3, \dots, t_J\}$   
 $\{w_1, w_2, w_3, \dots, w_I\}$ , LSA  
 $M_{ij}$  f  $w_i$   $t_j$ . S  
 (SVD) 10  
 $M = U\Sigma V^T$ . (1)  
 $MM^T$   $M^T M$   
 $\{w_1, w_2, \dots, w_I\}$ . T fi K  
 $\Lambda_K$  K  
 $MM^T$  K-  
 $w_i \rightarrow \Lambda_K(i)$  (2)  
 $\Lambda_K(i)$  i

**4. MOTIVATIONS OF USE OF SUBWORDS**

ASR  
 W  
 T f, LSA  
 H  
 f ASR  
 f C  
 f C

**4.1. Robustness to Flexibility in Word Segmentation**

AC  
 E  
 f

**Table 1. S f TDT2. E**

O	ASR	B
阿尔及利亚 (Algeria)	鲍尔 激励 要 (Bauer drive want)	a er ji li ya bao er ji li yao
奥尔布莱特 (Albright)	二步 莱特 (two step Wright)	ao er bu lai te er bu lai te
互联网 (internet)	互 连网 (mutual connection)	hu lian wang
赈济 (relieve)	震级 (quake magnitude)	zhen ji
过失 (defect)	国事 (national affair)	guo shi

ASR  
 T  
 F  
 北(N) 韩(K)  
 TDT2 M  
 T  
 H  
 ff

**4.2. Robustness to Speech Recognition Errors**

E C  
 1200  
 6500  
 400,  
 f homophones  
 C  
 T f C ASR  
 T I  
 TDT2  
 ASR F 阿尔及利亚(A) /鲍尔 激励 要/(B)  
 尔布赖特(A) /二步 莱特/(W). R  
 H  
 F 阿尔及利亚 /鲍尔 激励 要/,  
 bu lai te 赖特 奥  
 尔布赖特 /二步 莱特/

**4.3. Robustness to OOV Words**

T fi C  
 H  
 C  
 C ASR  
 BN  
 f  
 N  
 OOV  
 10% f BN  
 A  
 OOV  
 ff  
 f  
 f

Table 2. S f OOV v. f TDT2. E

C	B
OOV v. : 王有才 (a Chinese name)	wang you cai
ASR	当有财 (when have money)
	干油菜 (king rape)
	邦有才 (national friendship talent)
OOV v. : 莱温斯基 (Lewinsky)	lai wen si ji
ASR	来文斯基 (come article this base)
	来问司机 (come ask driver)
	来的司机 (show-up driver)
OOV v. : 科索沃 (Kosovo)	ke suo wo
ASR	克租同 (gram motherland)
	客座我 (guest me)

... F, f ... OOV C ...  
 ... A ...  
 ... T 2 v. f ...  
 OOV v. f TDT2. F ... OOV f  
 莱温斯基 (L v. ) ff  
 ... /米问司机, /米文斯基 /米  
 的司机, v. f M  
 v. f OOV v. ASR ( ).

5. WORD AND SUBWORD LSA IN TEXTTILING

5.1. TextTiling-based Story Segmentation

T T T  
 7: A  
 fi . A  
 ASR TDT2  
 f v. S 5.3. T f v.  
 I f fi  
 g lexical score:

$$(g) = \cos(\mathbf{v}_s, \mathbf{v}_{s+1}) = \frac{\sum_{i=1}^I v_{s,i} v_{s+1,i}}{\sqrt{\sum_{i=1}^I v_{s,i}^2 \times \sum_{i=1}^I v_{s+1,i}^2}} \quad (3)$$

v. v<sub>s</sub>, v<sub>s+1</sub> f f v.  
 s s+1 f v<sub>s</sub>, ... f f w<sub>i</sub> s. I v.  
 fi fi f (T). T ASR f  
 f f ;2) f v. v.  
 7. S  
 L v. Δ {T, T+Δ, T+2Δ...}  
 T T depth score Δ ≤ T.  
 fi D

$$(u) = (p_l) - (u) + (p_r) - (u) \quad (4)$$

v. u, , p<sub>l</sub> p<sub>r</sub> f  
 u, . T  
 F , fi  
 f fi v. v.  
 - fi θ

5.2. Applying LSA to TextTiling

I , LSA f f  
 12. A S 5.1, T T . T LSA f T  
 T f f . A ASR  
 ASR  
 ( 0.8  
 LSA f f  
 T LSA  
 T f f v. v. .1) S  
 1.  
 2) C. fi v. v.  
 v. f v. f- f  
 T v. f- f

W LSA ( S 3)  
 SVDLIBC  
 LSA f Λ<sub>K</sub>(i) f (2)  
 T T T T f W Λ<sub>K</sub>(i)  
 S 5.1) K- ( fi

$$\hat{\mathbf{v}}_s = \sum_{i=1}^I v_{i,s} \times \Lambda_K(i), \quad (5)$$

$$(g) = \cos(\hat{\mathbf{v}}_s, \hat{\mathbf{v}}_{s+1}) = \frac{\sum_{i=1}^K \hat{v}_{s,i} \hat{v}_{s+1,i}}{\sqrt{\sum_{i=1}^K \hat{v}_{s,i}^2 \times \sum_{i=1}^K \hat{v}_{s+1,i}^2}}, \quad (6)$$

v. v̂<sub>s,i</sub> f v̂<sub>s</sub>. S  
 fi w<sub>i</sub> f E . (6). S fi  
 v. -LSA- T T

5.3. Subword-LSA-based TextTiling

W f LSA ff C v. F  
 f v. {w<sub>1</sub>w<sub>2</sub>w<sub>3</sub>...w<sub>Q</sub>},  
 fi f ( )  
 ) ( ), ..., {c<sub>1</sub>c<sub>2</sub>c<sub>3</sub>...c<sub>L</sub>}. T  
 v. f  
 : {c<sub>1</sub>c<sub>2</sub> c<sub>2</sub>c<sub>3</sub> c<sub>3</sub>c<sub>4</sub>...c<sub>L-1</sub>c<sub>L</sub>}, (7)  
 : {c<sub>1</sub>c<sub>2</sub>c<sub>3</sub> c<sub>2</sub>c<sub>3</sub>c<sub>4</sub> c<sub>3</sub>c<sub>4</sub>c<sub>5</sub>...c<sub>L-2</sub>c<sub>L-1</sub>c<sub>L</sub>}, (8)

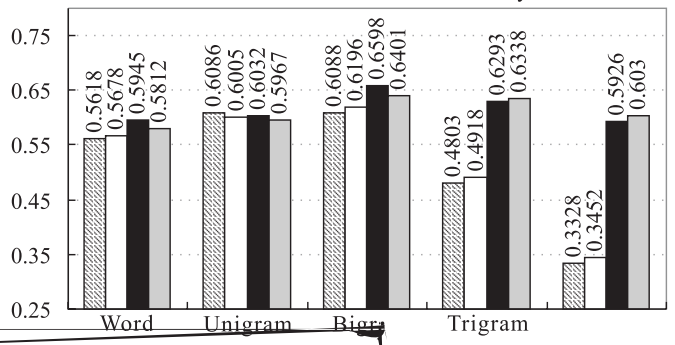
A S 5.2, v. LSA  
 f v - f  
 , v.  $w_i$   
 . L E . (6)

### 6. EXPERIMENTS

W (1)  
 v. T T 6, (2)  
 v. T T 7, (3) v. -LSA- T T (4)  
 v. -LSA- T T E v.  
 fi f

F1- f P  
 $fT = 50 \quad \Delta = 20$   
 v. f T T T  
 T  $\Delta$  v. fi f  
 $\theta$  LSA  $Kf$  v. v.  
 2. E  
 f F1- v. F .1.  
 R v. LSA T  
 T f C  
 BN LSA f  
 v. T T v.  
 LSA.B -LSA- T T f  
 C -LSA- F1- f 0.6598 v.  
 f 17.4% v.  
 T T ( f v. ) 6.5%  
 - - - T T T f f.  
 A (2957 f 395 f ) f )  
 T (LSA f ) f /  
 LSA v. f f v.  
 C - f v.

▨ character w/o LSA □ syllable w/o LSA  
 ■ character w/ LSA ▩ syllable w/ LSA



T T  
 v. T  
 (LSA)  
 LSA  
 LSA,

0.6598 v.