



- 이 그림은... 그...  
옛날 그림인데... 멋져요!

# CROSS – MODAL BIDIRECTIONAL TRANSLATION VIA REINFORCEMENT LEARNING

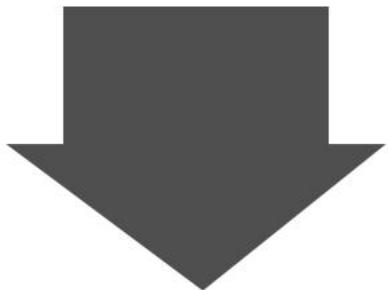
2020. 06. 10.  
Moon Sungwon

# INTRODUCTION

---

TEXT

이 영화 진짜 재미 없다



**Class: 부정 (99%)**

Image



**Class: 강아지 (87%)**

# INTRODUCTION

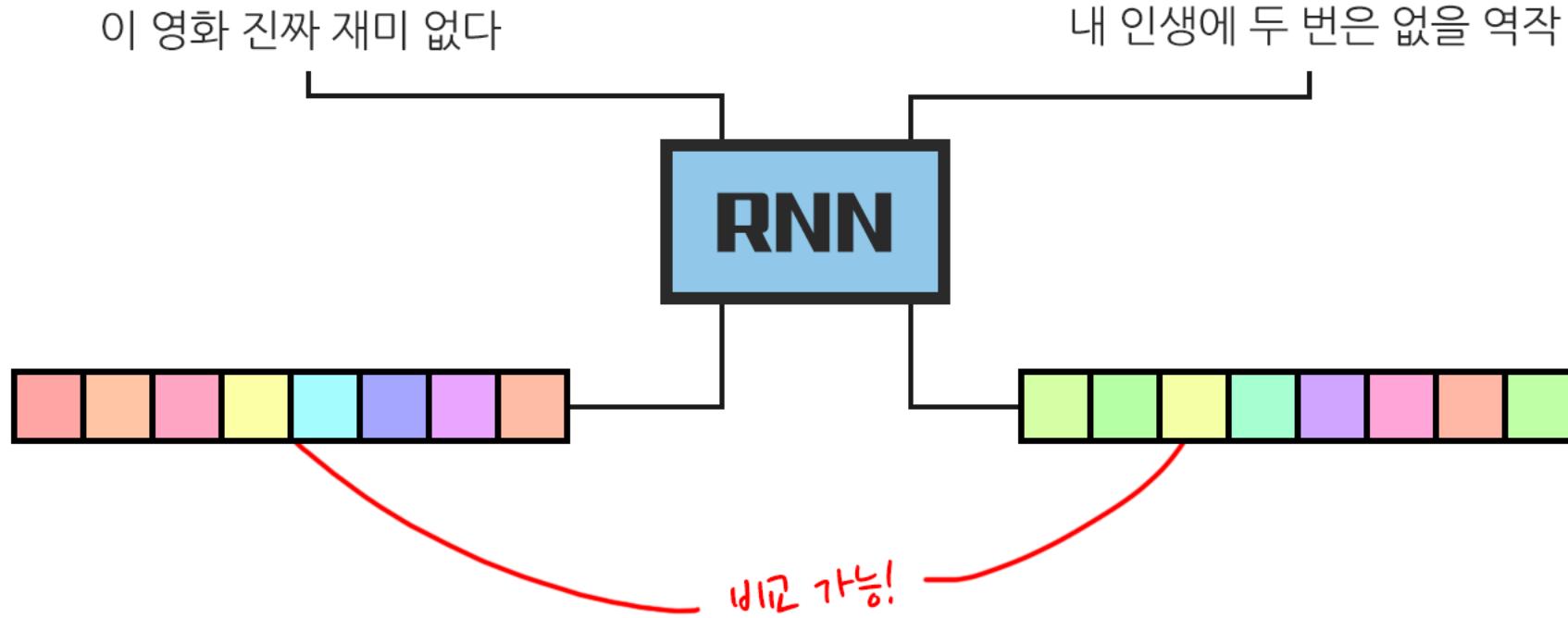
---

이 영화 진짜 재미 없다



# INTRODUCTION

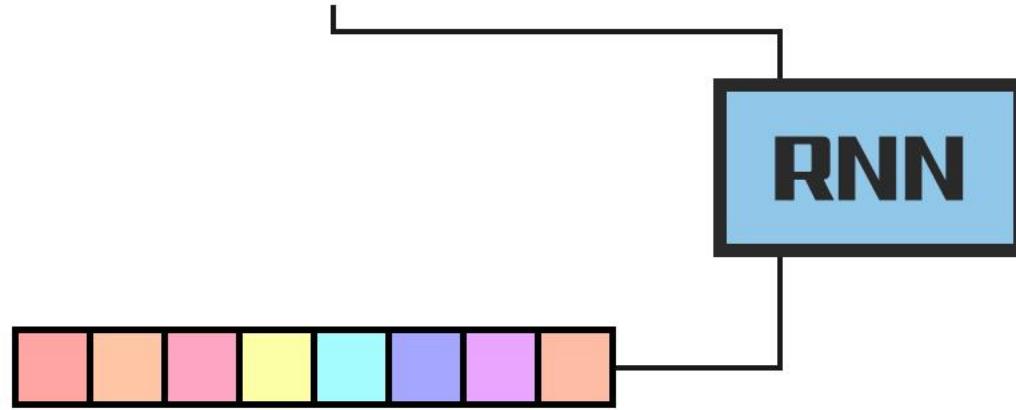
---



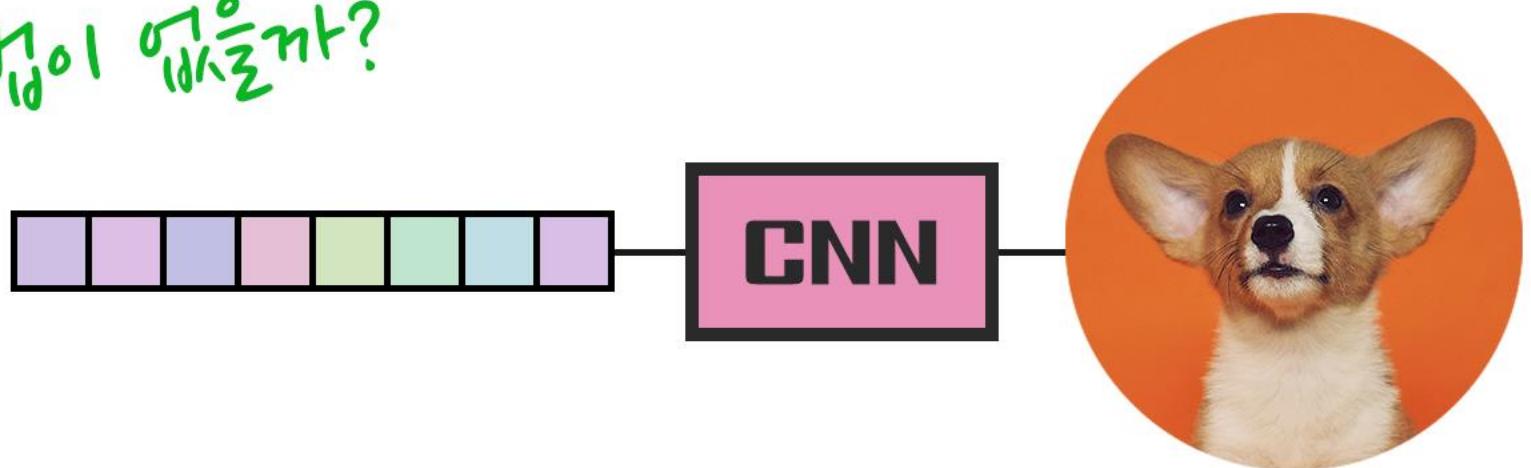
# INTRODUCTION

---

이 영화 진짜 재미 없다



방법이 없을까?



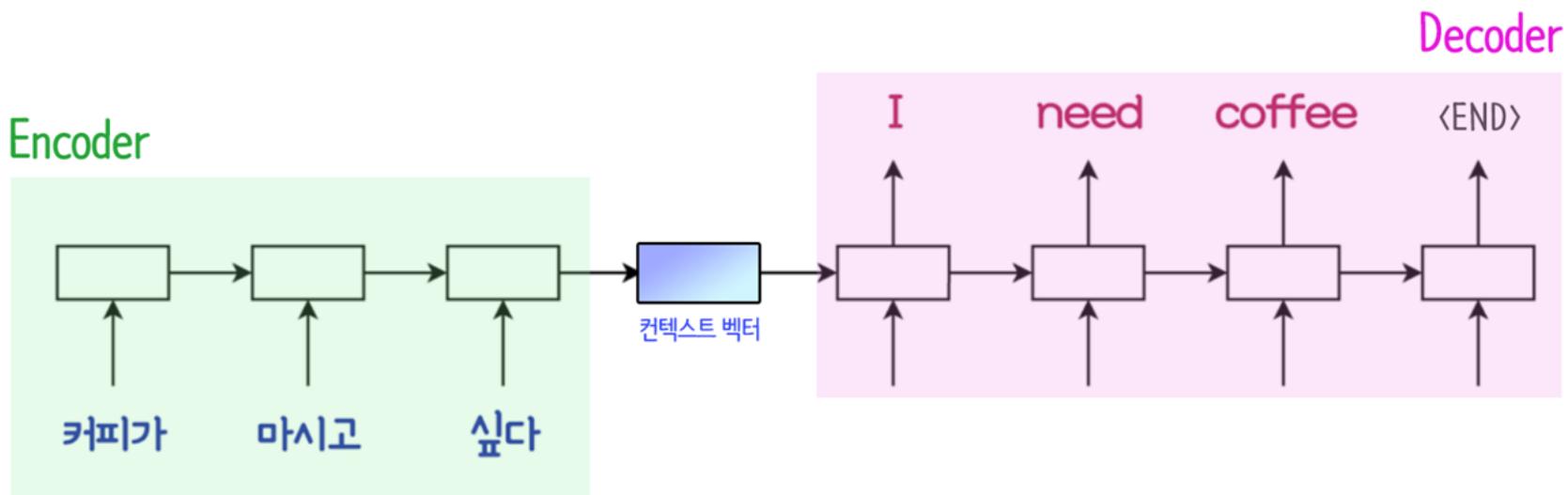
# INTRODUCTION

---

“ **Cross-modal Translation** ”  
= Seq2Seq

# INTRODUCTION

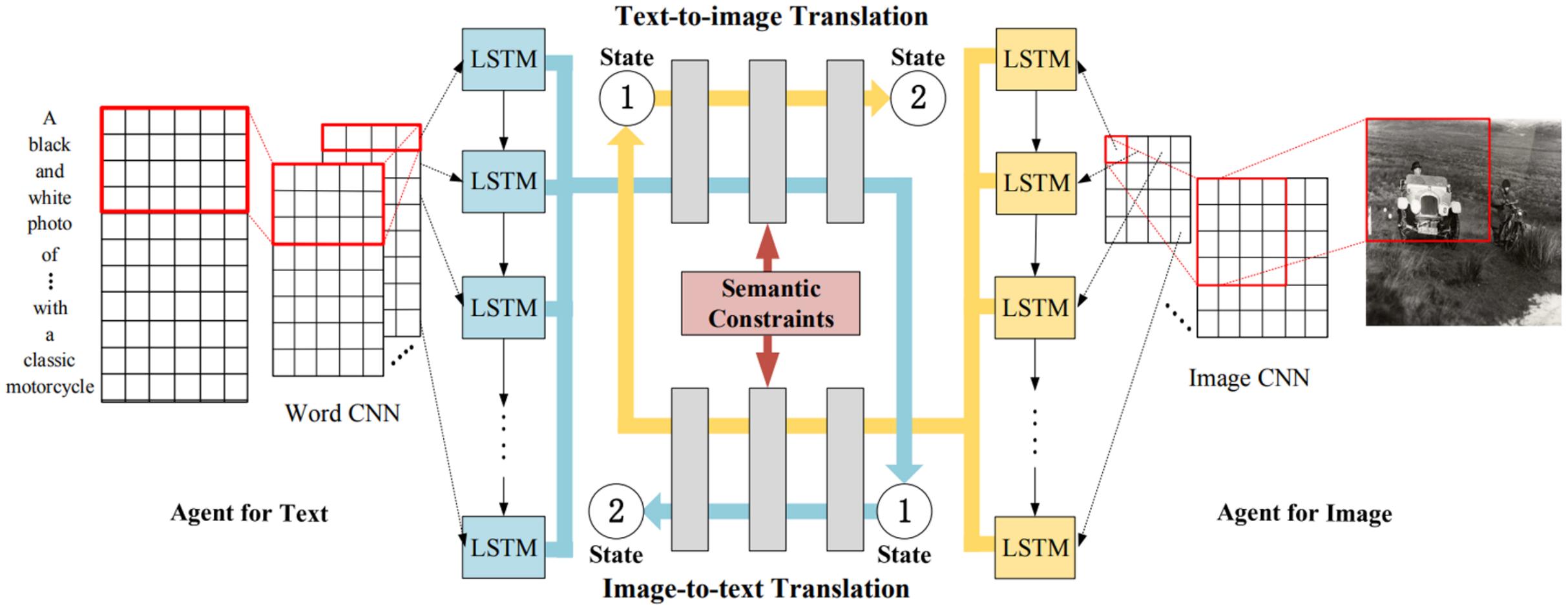
기존의 방법이 대체 어땠길래?



▲ Sequence-to-Sequence (Seq2seq)

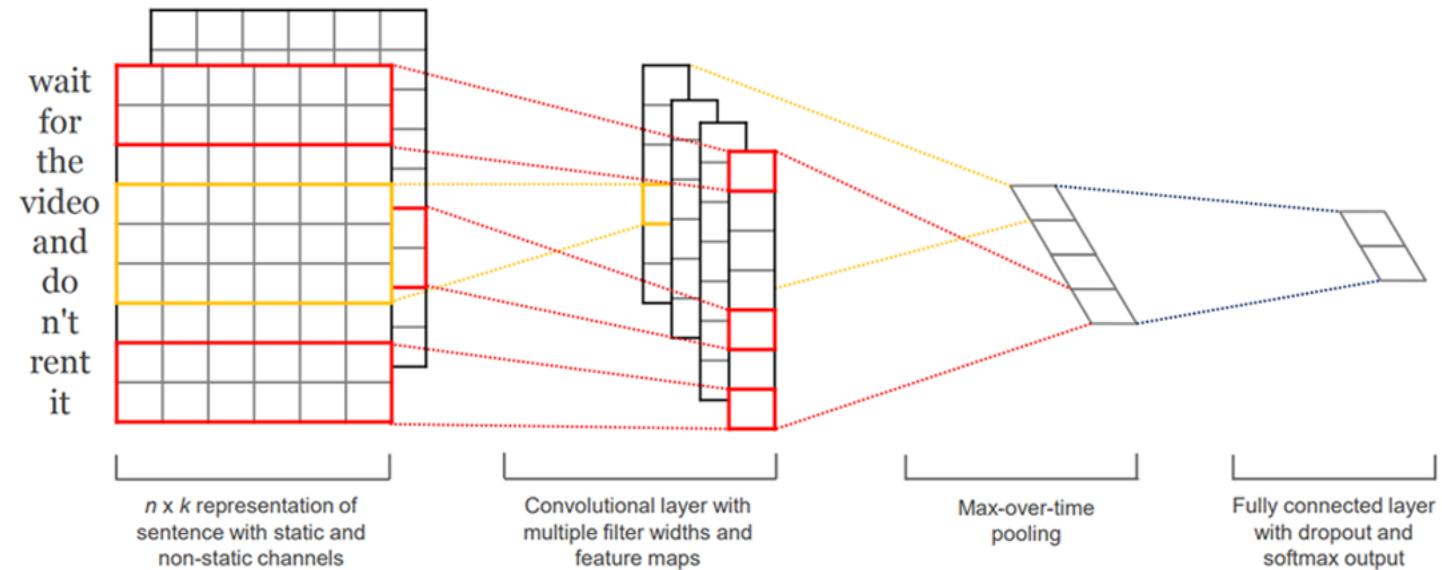
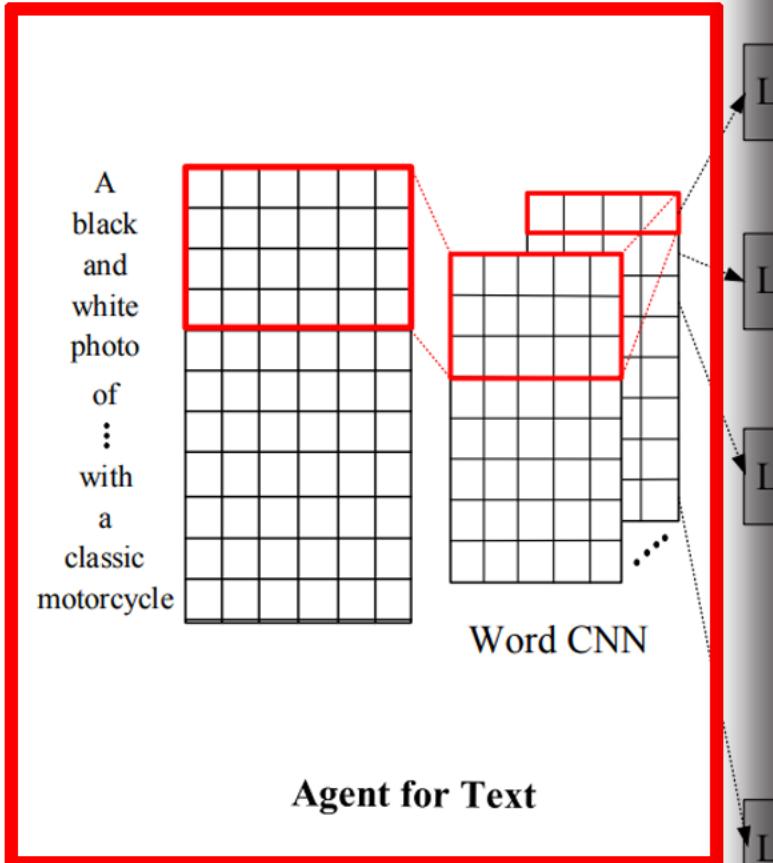
# MODEL

---



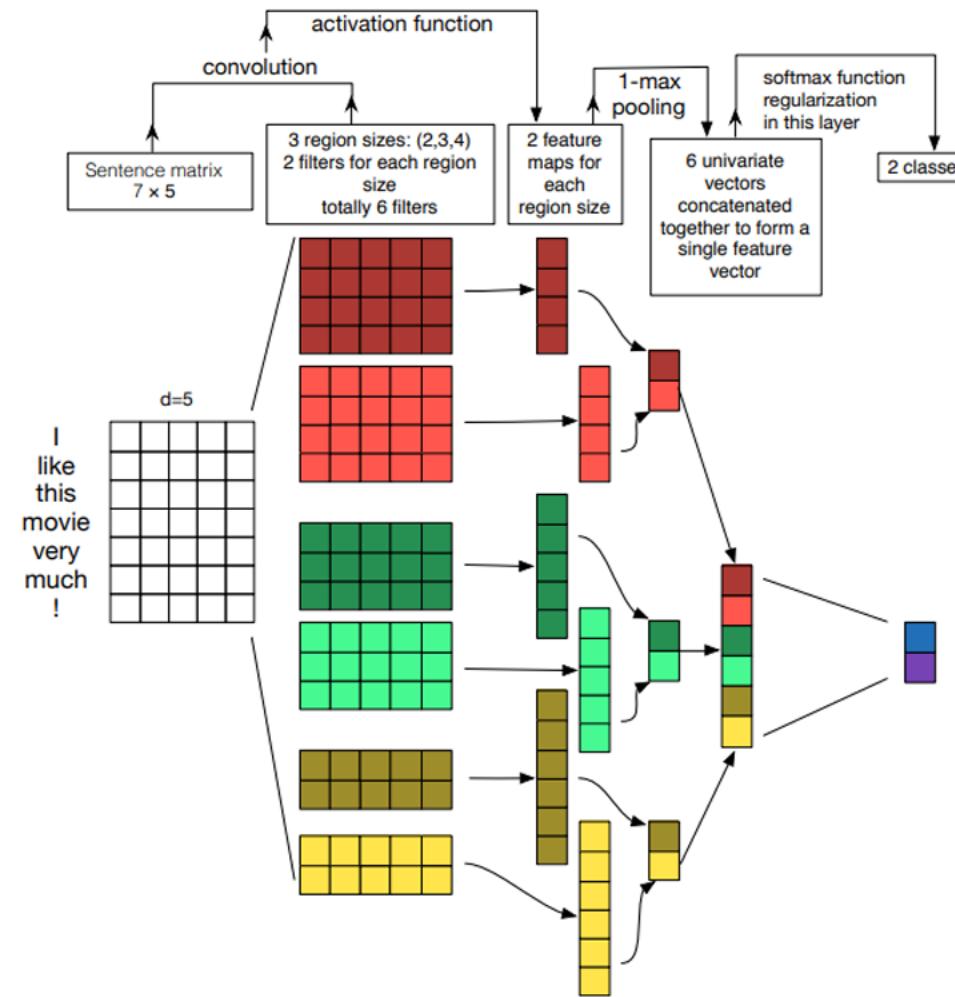
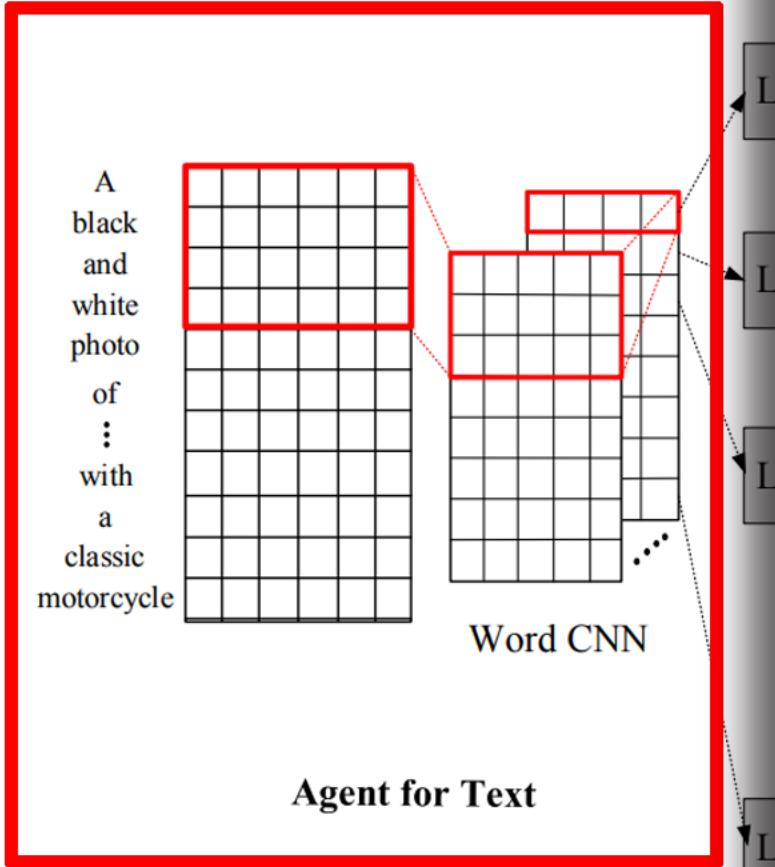
# MODEL

---



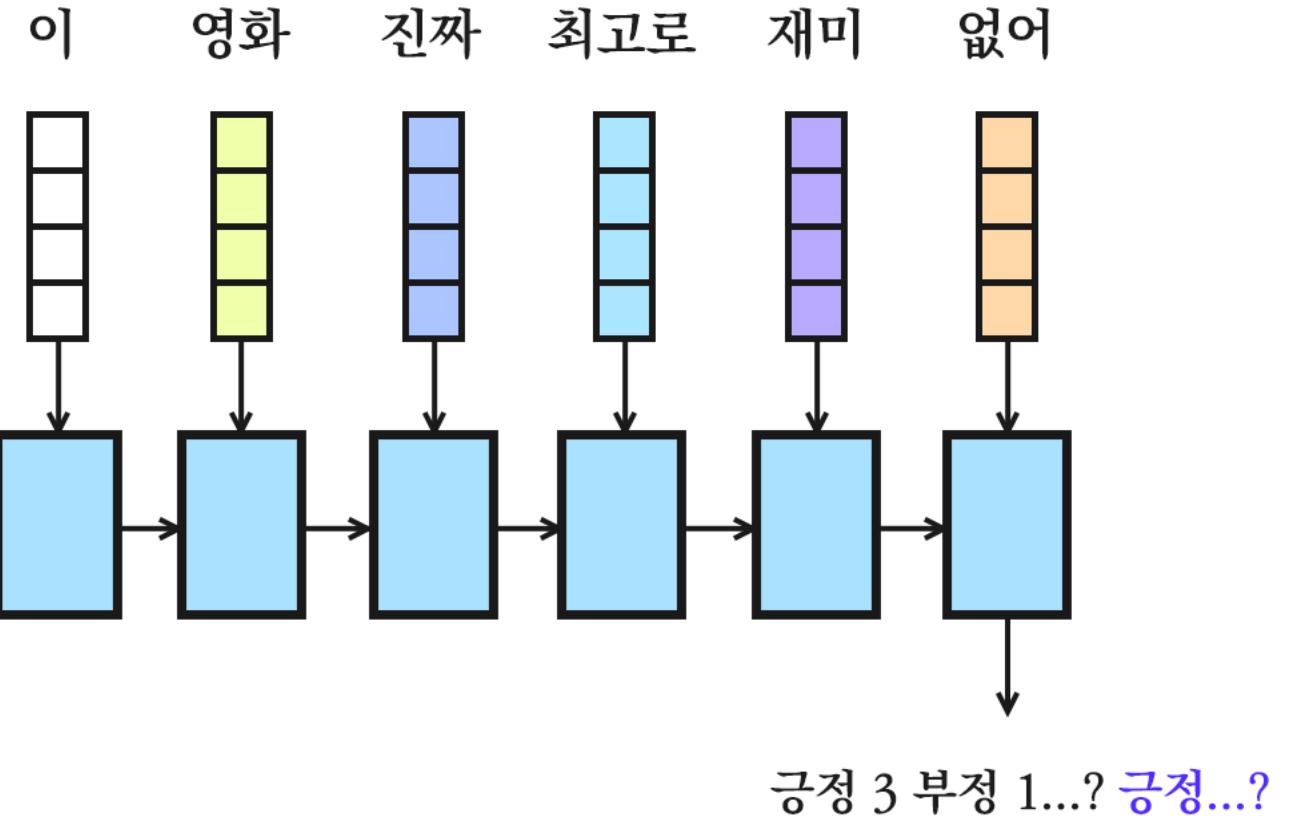
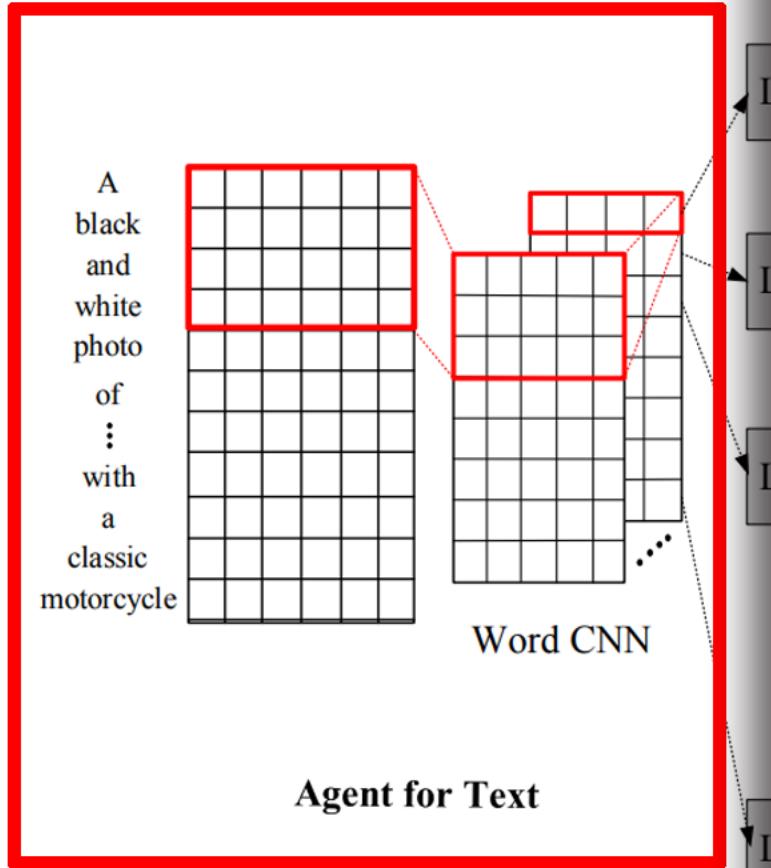
Convolutional Neural Networks for Sentence Classification (2014)

# MODEL



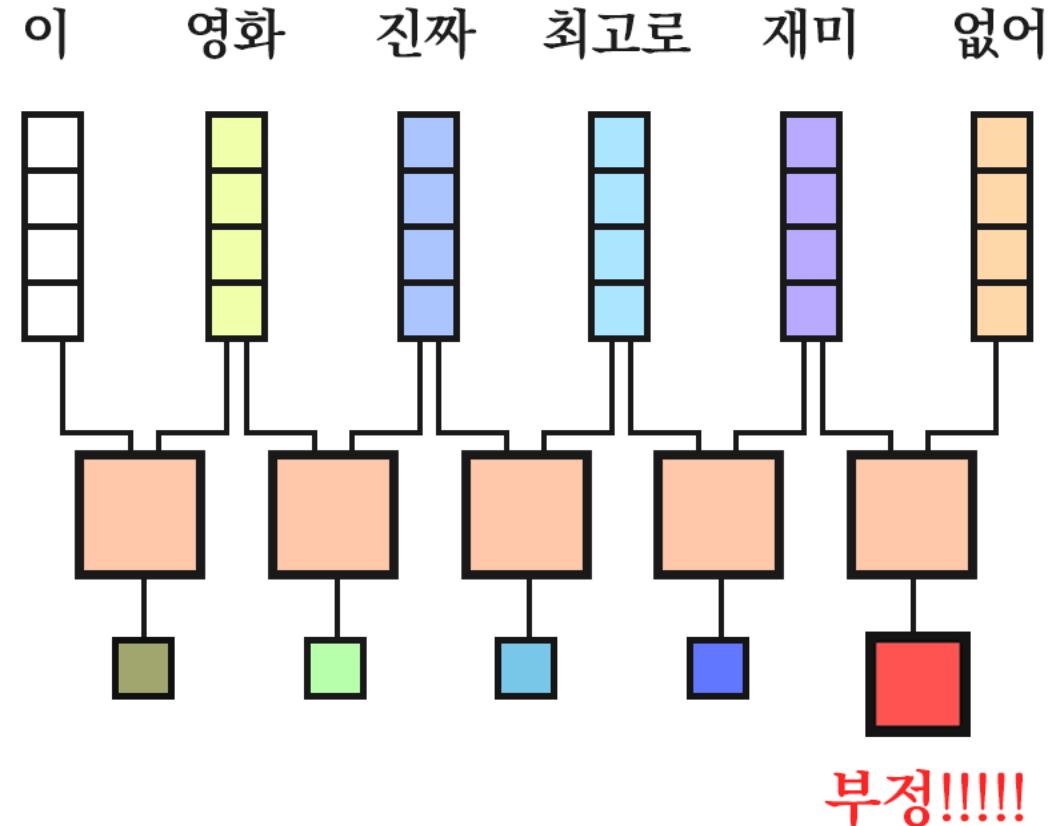
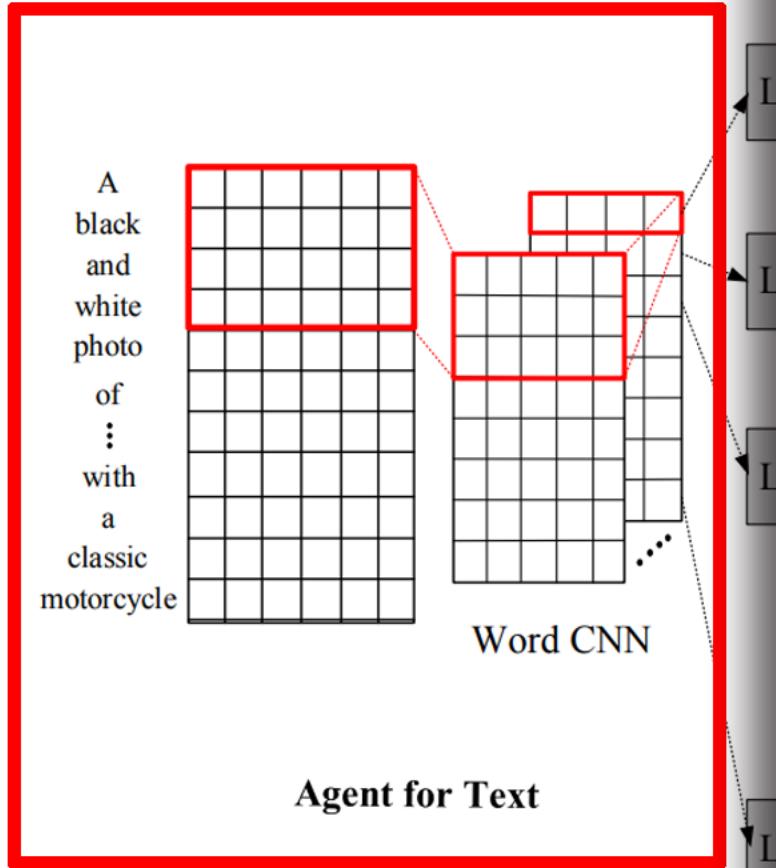
A Sensitivity Analysis of  
(and Practitioners' Guide to)  
Convolutional Neural Networks  
for Sentence Classification  
(2015)

# MODEL



## Unigram Model

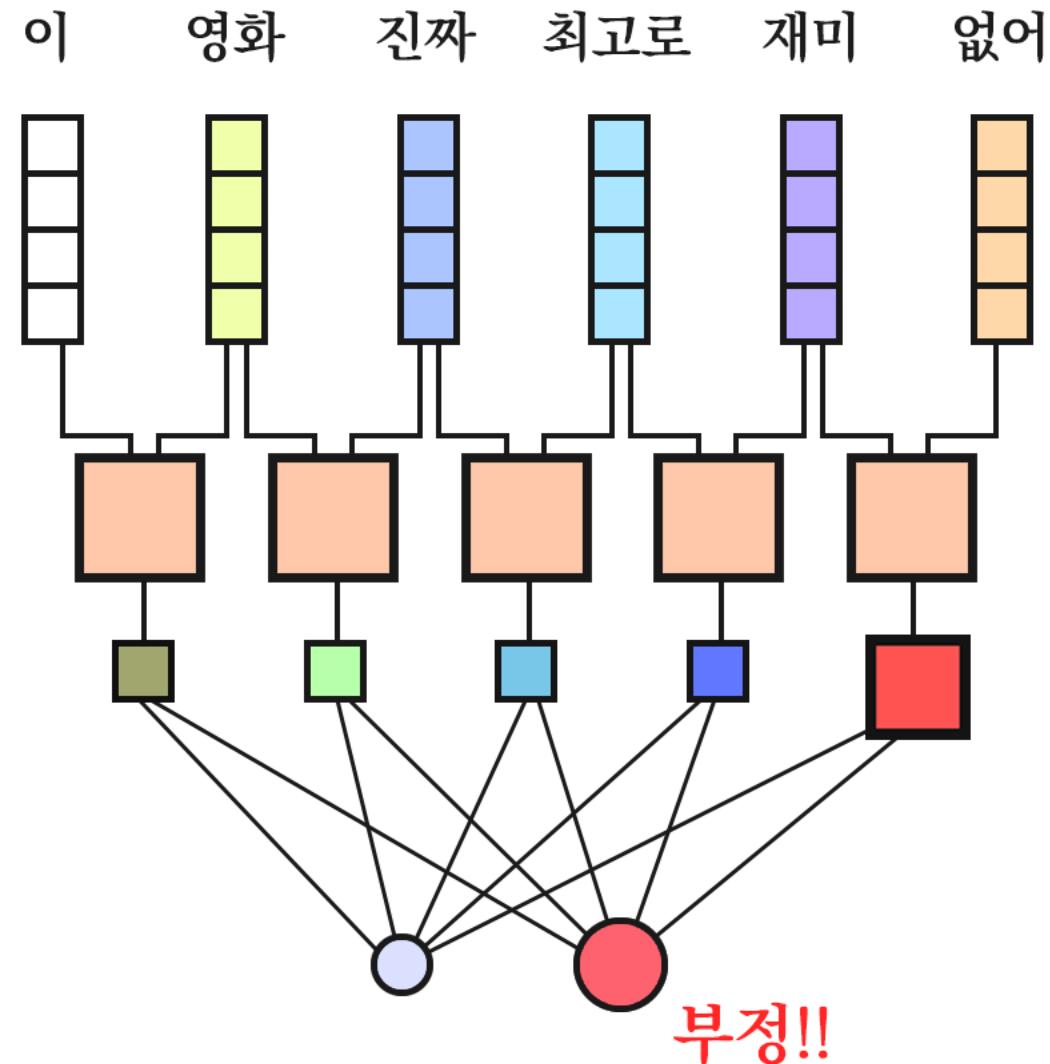
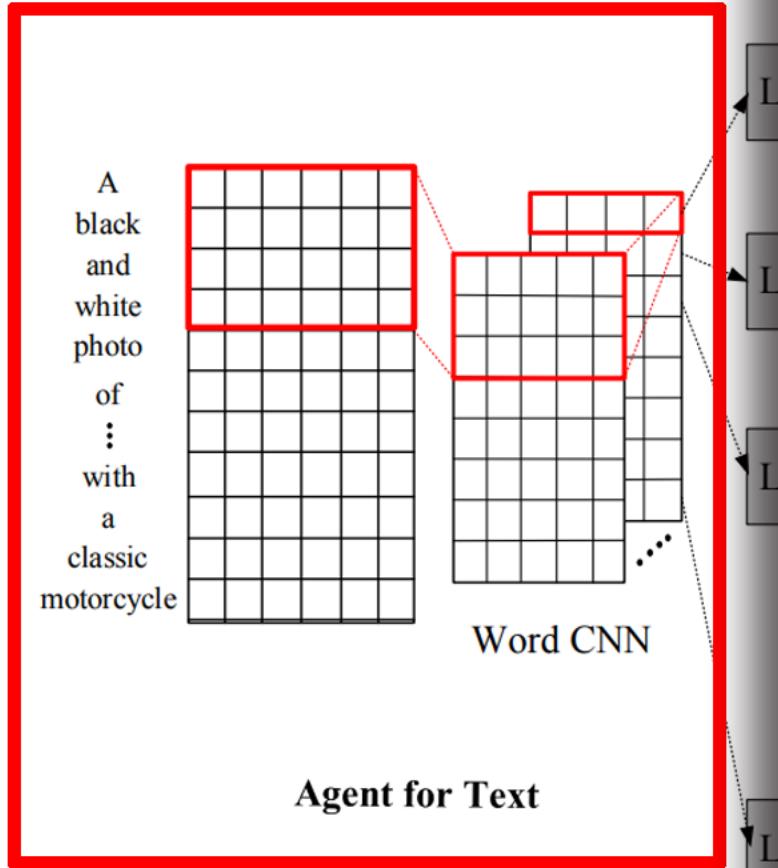
# MODEL



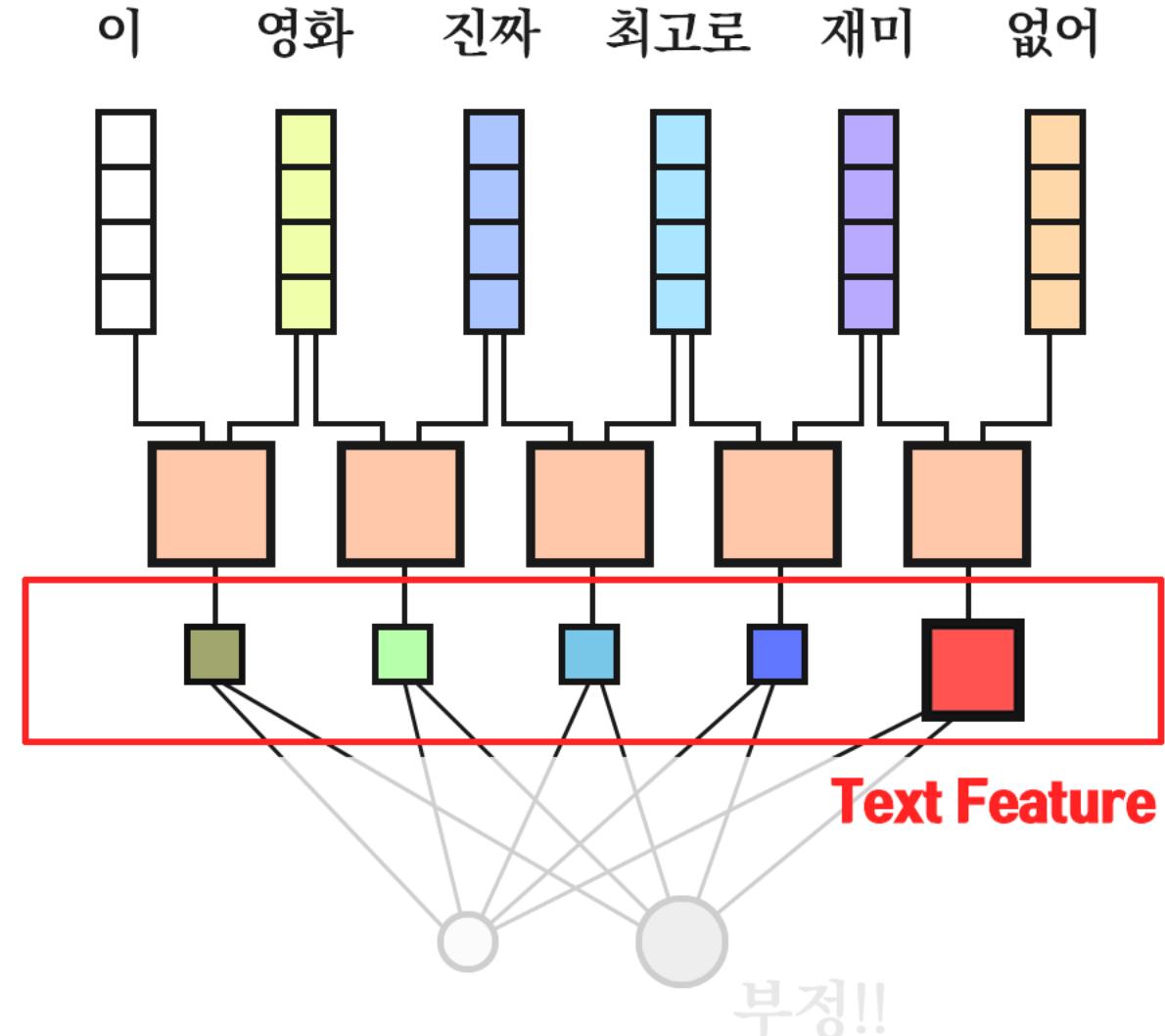
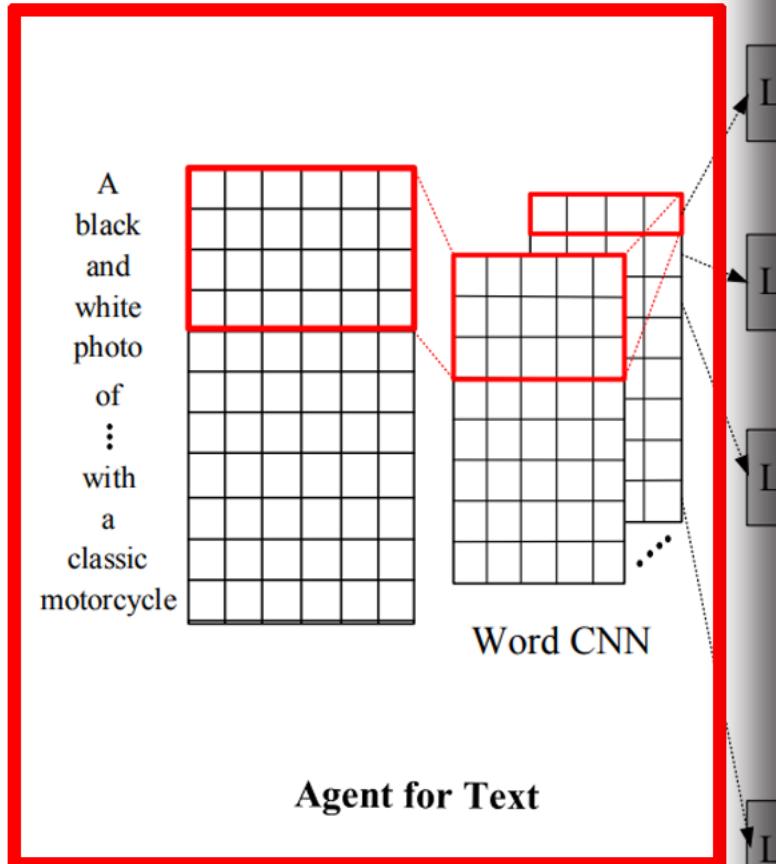
## N-Gram Model

: 문맥을 파악하는 데에 유리!

# MODEL

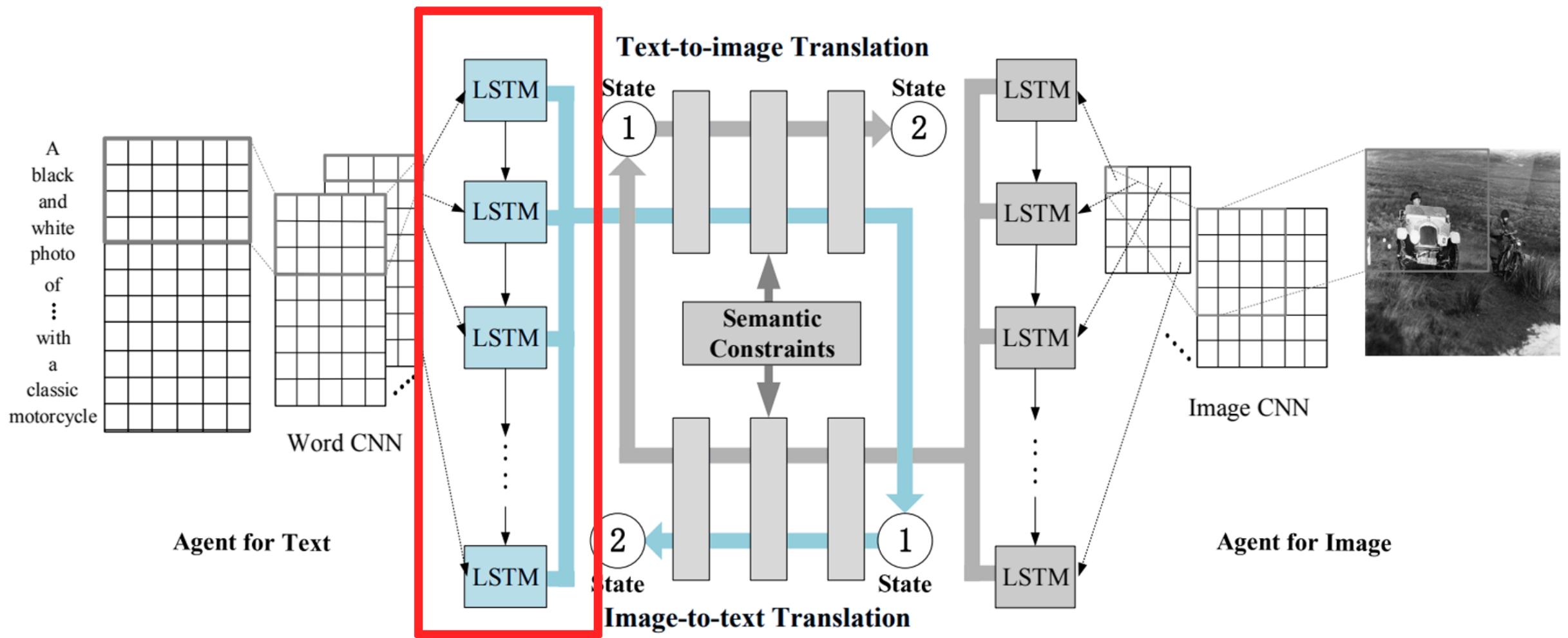


# MODEL

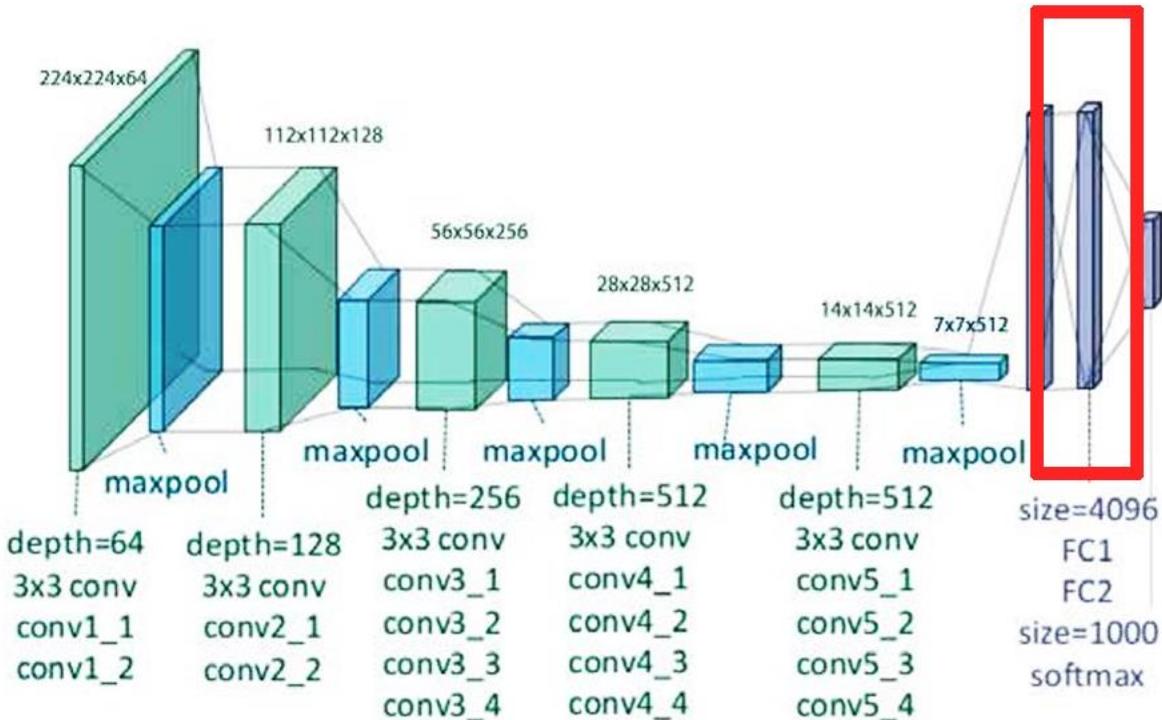


# MODEL

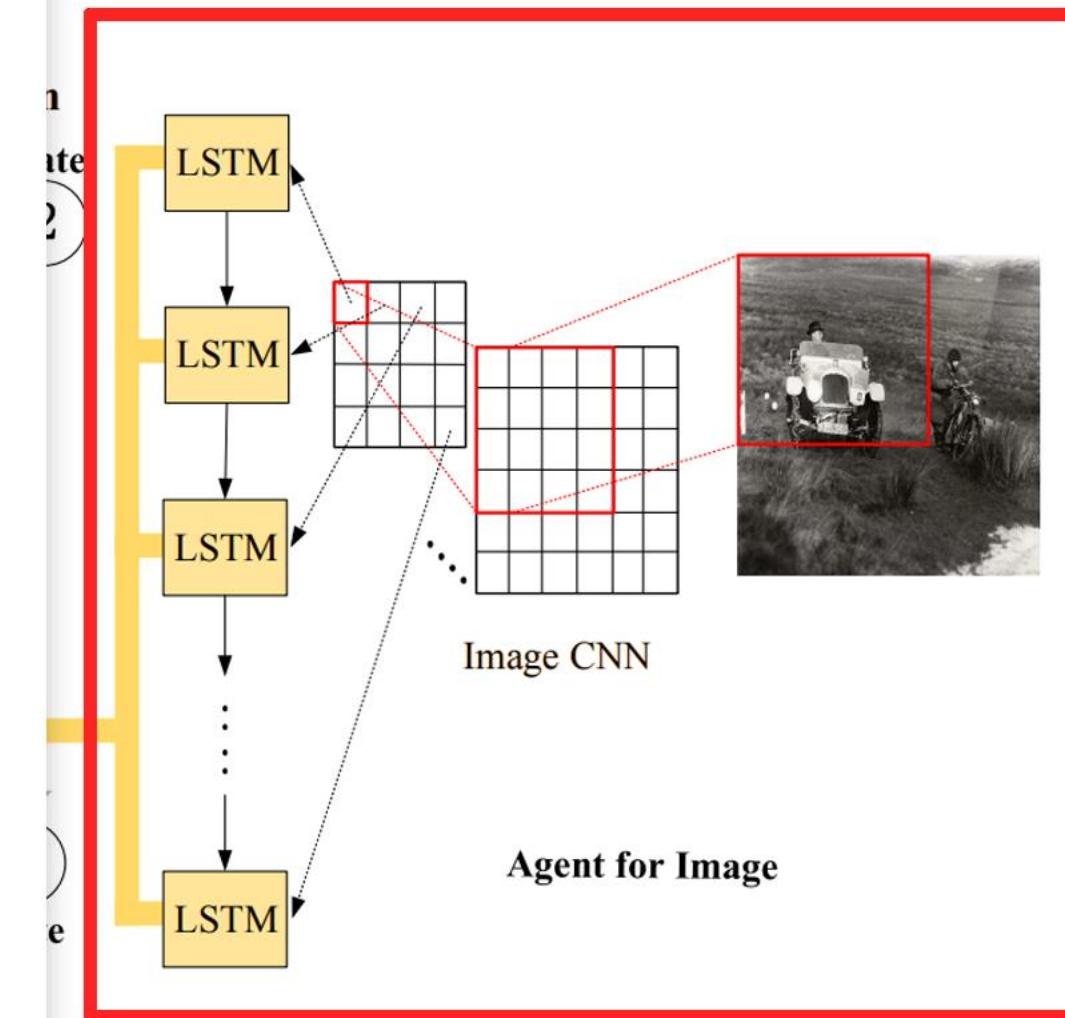
---



# MODEL

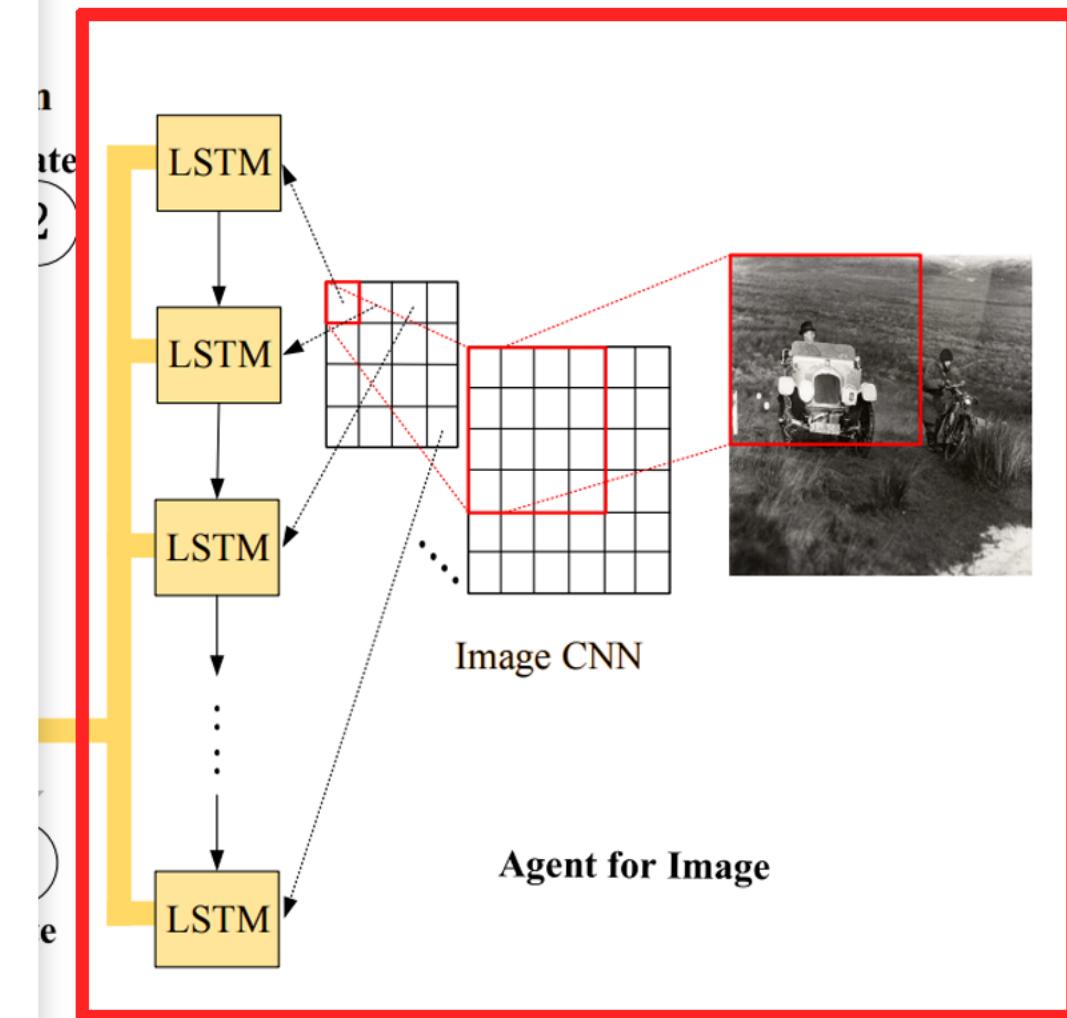


VERY DEEP CONVOLUTIONAL NETWORKS  
FOR LARGE-SCALE IMAGE RECOGNITION  
(2014)

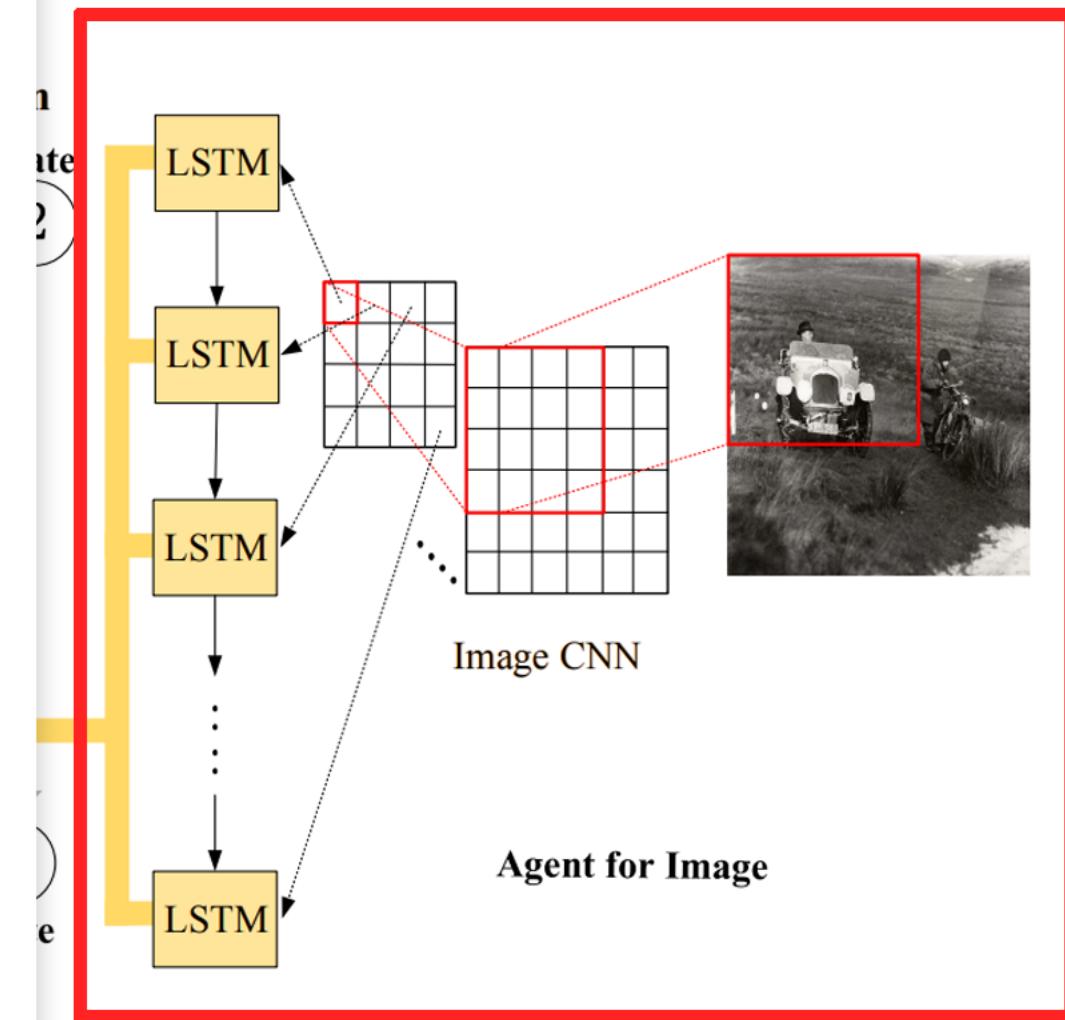


# MODEL

VGG-19 Input size



# MODEL

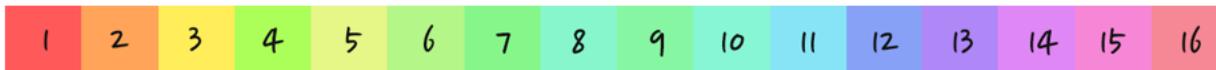


# MODEL

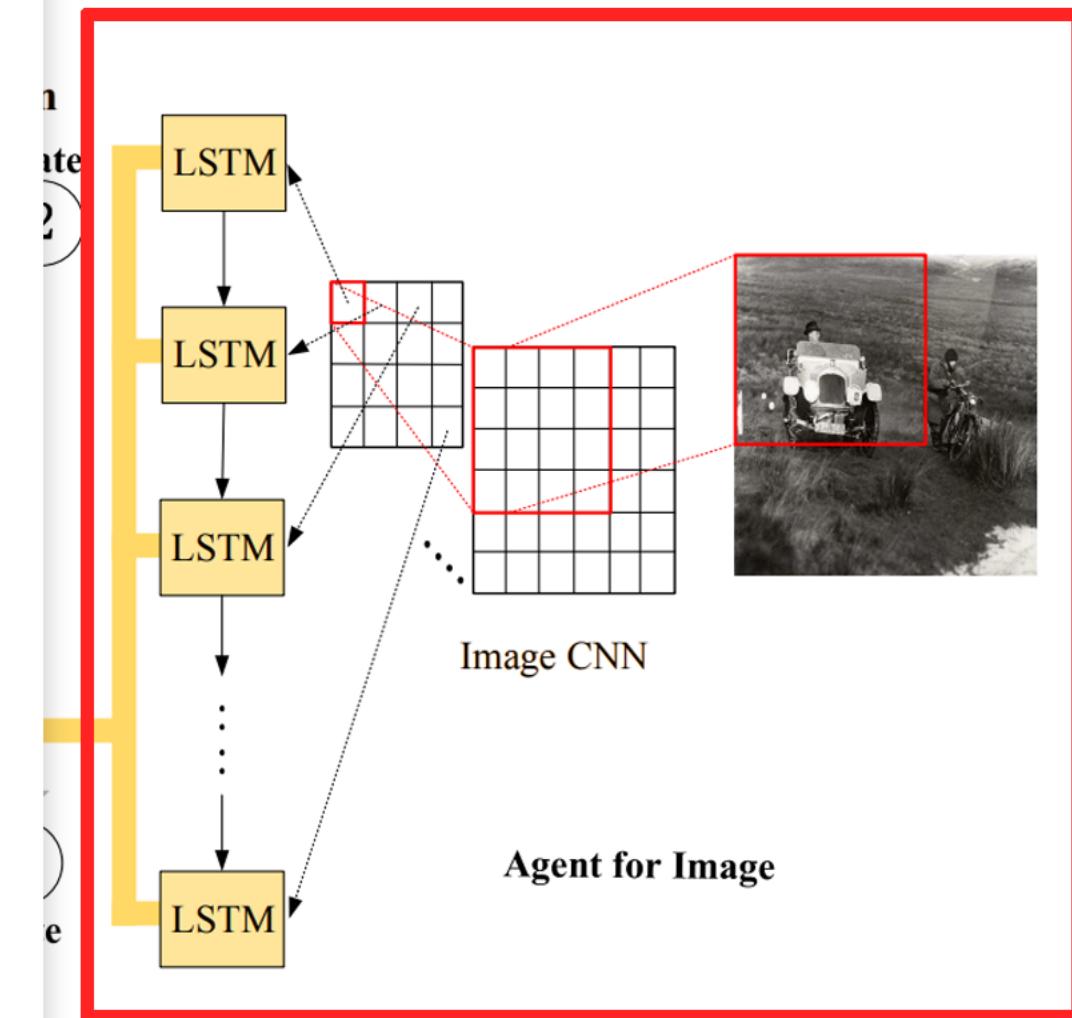


Feature Map

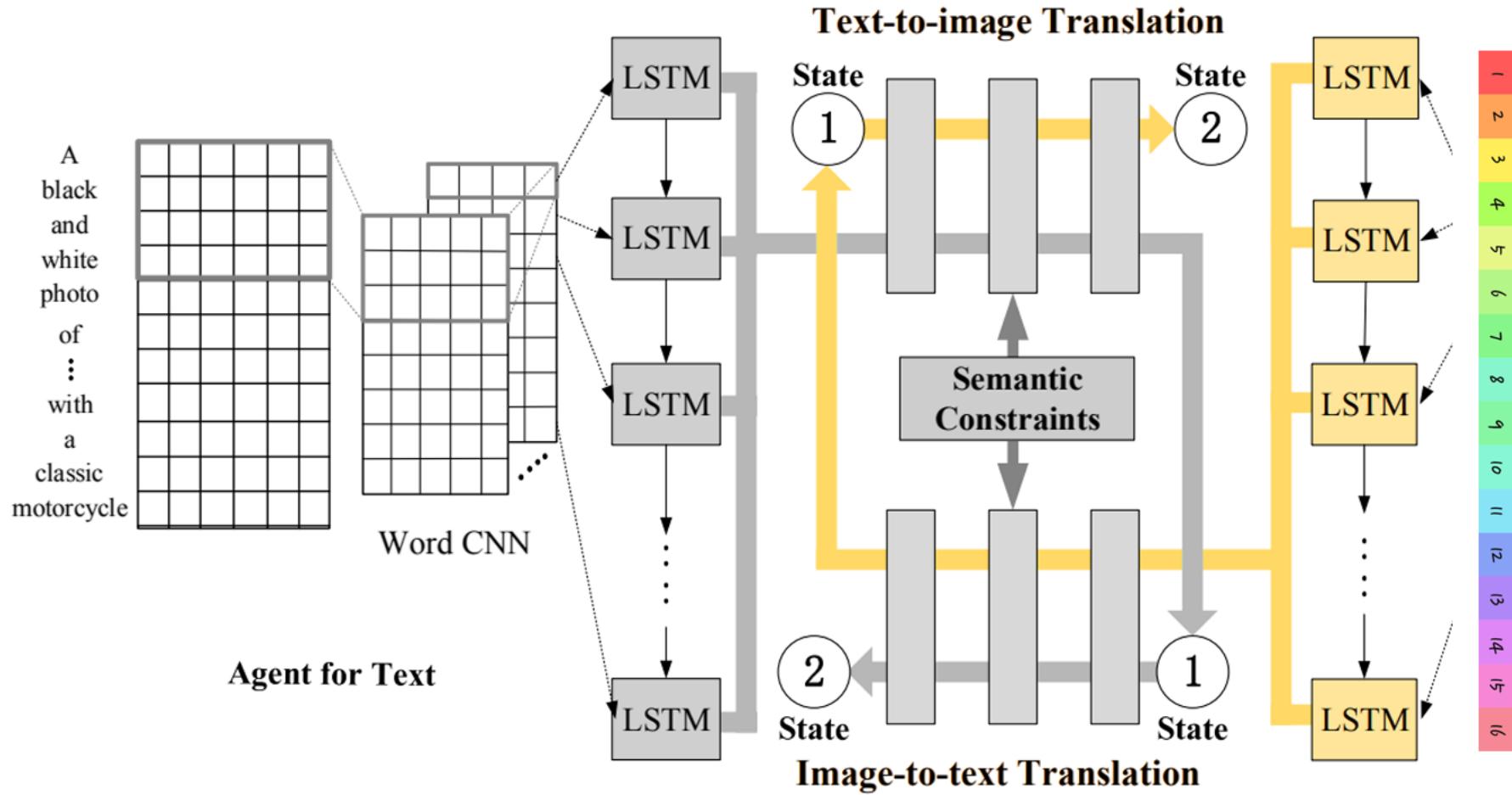
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



시선의 이동

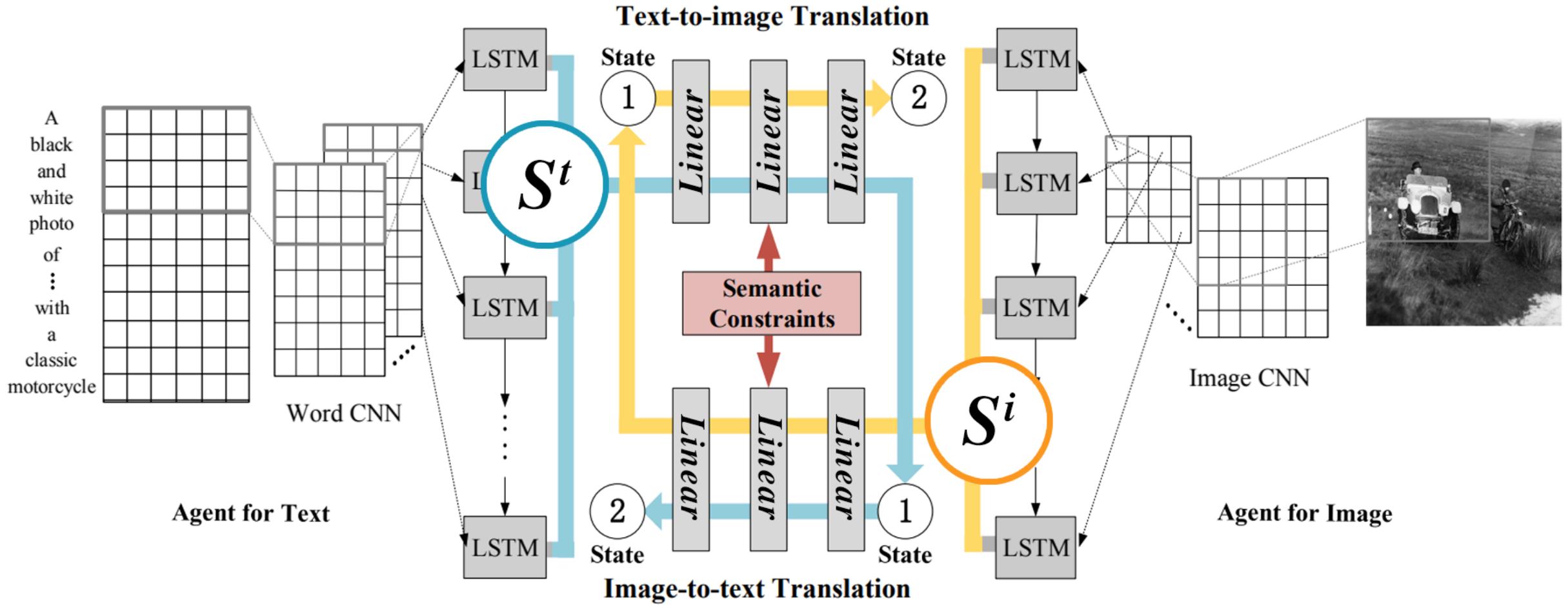


# MODEL



# MODEL

---



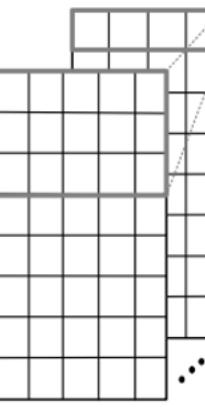
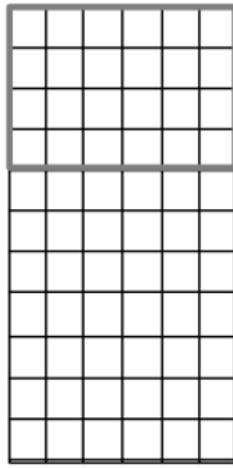
# MODEL

*State:  $S^t$*

*Action: Linear*

$\rightarrow S_{mid}^i \leftrightarrow S^i \rightarrow \text{Reward}^{\text{inter}}$

A black and white photo of : with a classic motorcycle



Word CNN

Agent for Text

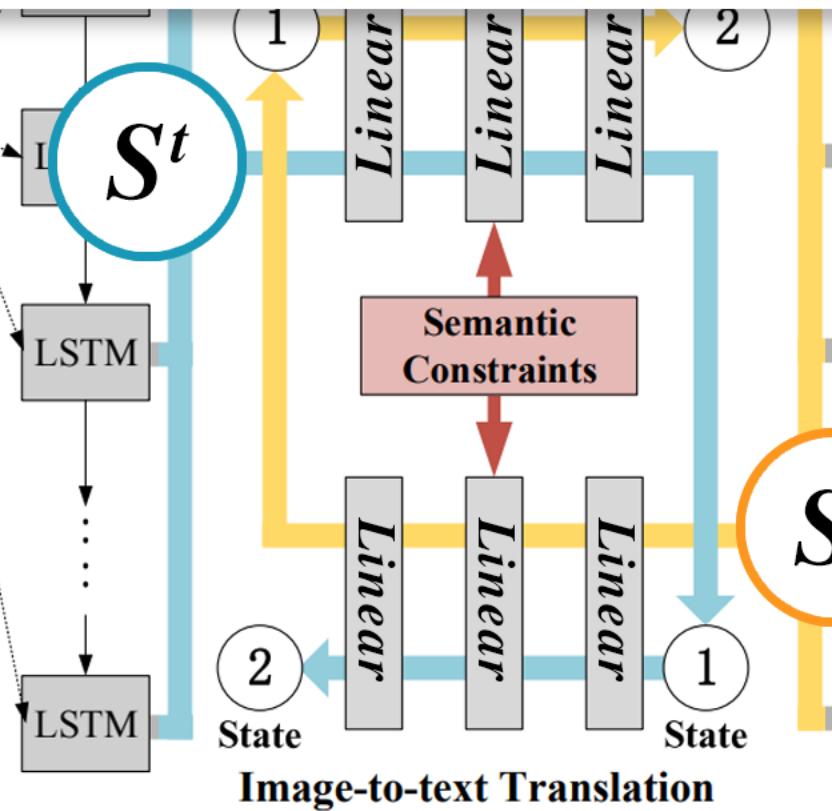


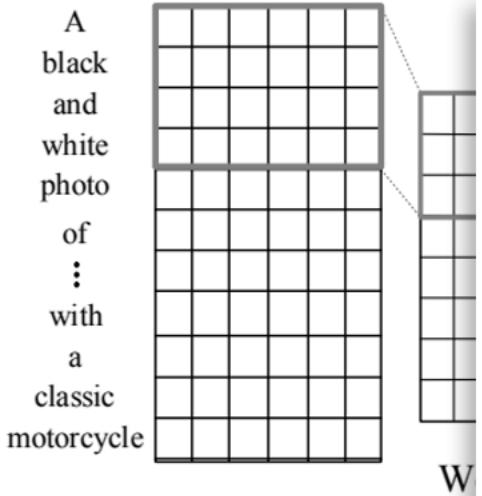
Image CNN

Agent for Image



# MODEL

**State:**  $S_{mid}^i$   
**Action:** Linear  
→  $S_{ori}^t$  ↔  $S^t$  → **Reward<sup>intra</sup>**



Predict

$$r_p^{inter} = \log(\text{norm}\left(\frac{s_{mid,p}^i \cdot s_p^t}{\|s_{mid,p}^i\|_2 \|s_p^t\|_2}\right)) \quad \text{Real}$$
$$r_p^{intra} = \log(\text{norm}\left(\frac{s_p^i \cdot s_{ori,p}^i}{\|s_p^i\|_2 \|s_{ori,p}^i\|_2}\right))$$



Agent for Te

nage



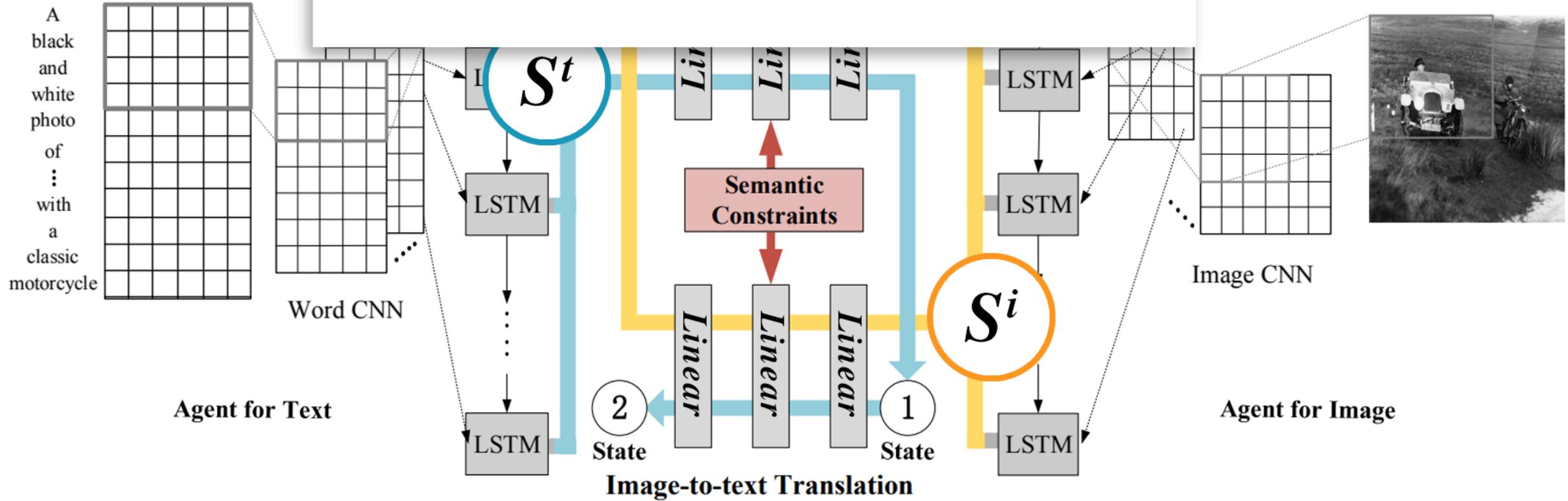
# MODEL

*Text* → *Image* → *Text*

$R^{inter}$

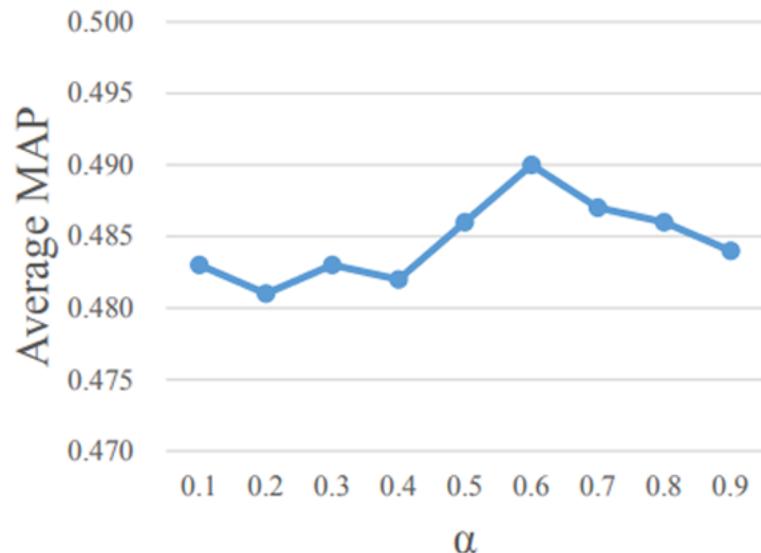
$R^{intra}$

$$r_p = \alpha r_p^{inter} + (1 - \alpha) r_p^{intra}$$

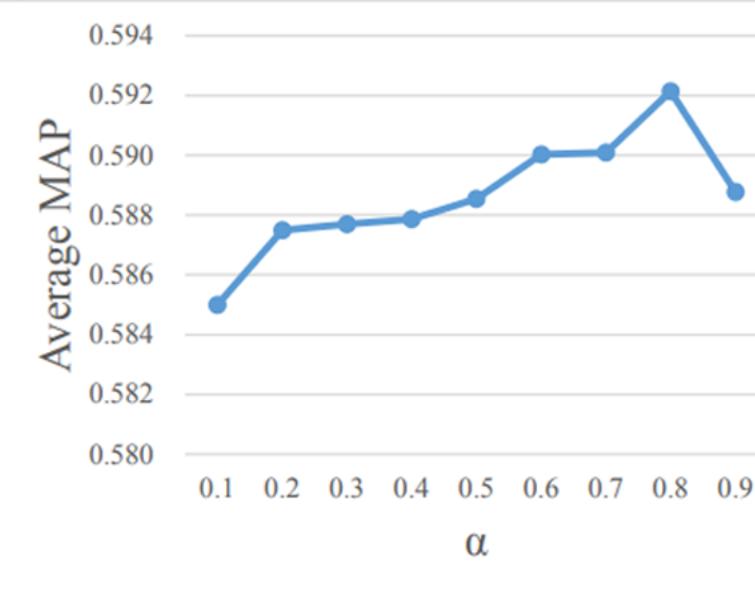


# MODEL

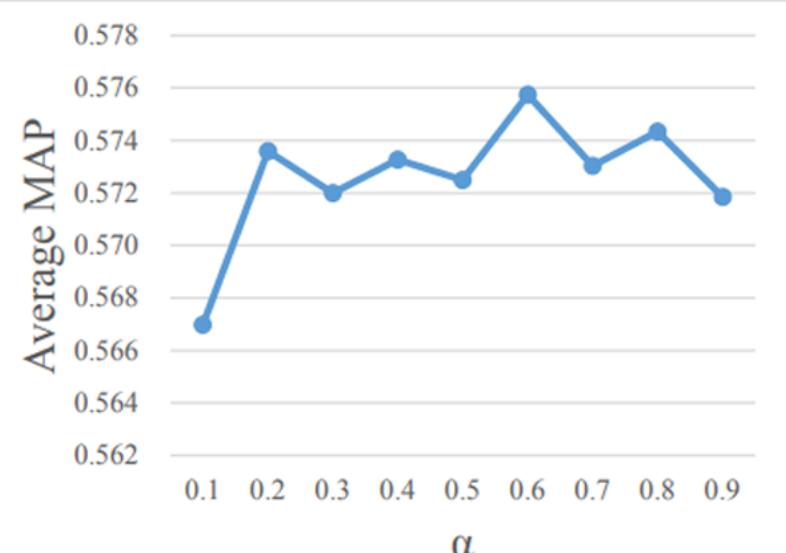
---



(a) Wikipedia dataset



(b) Pascal Sentence dataset

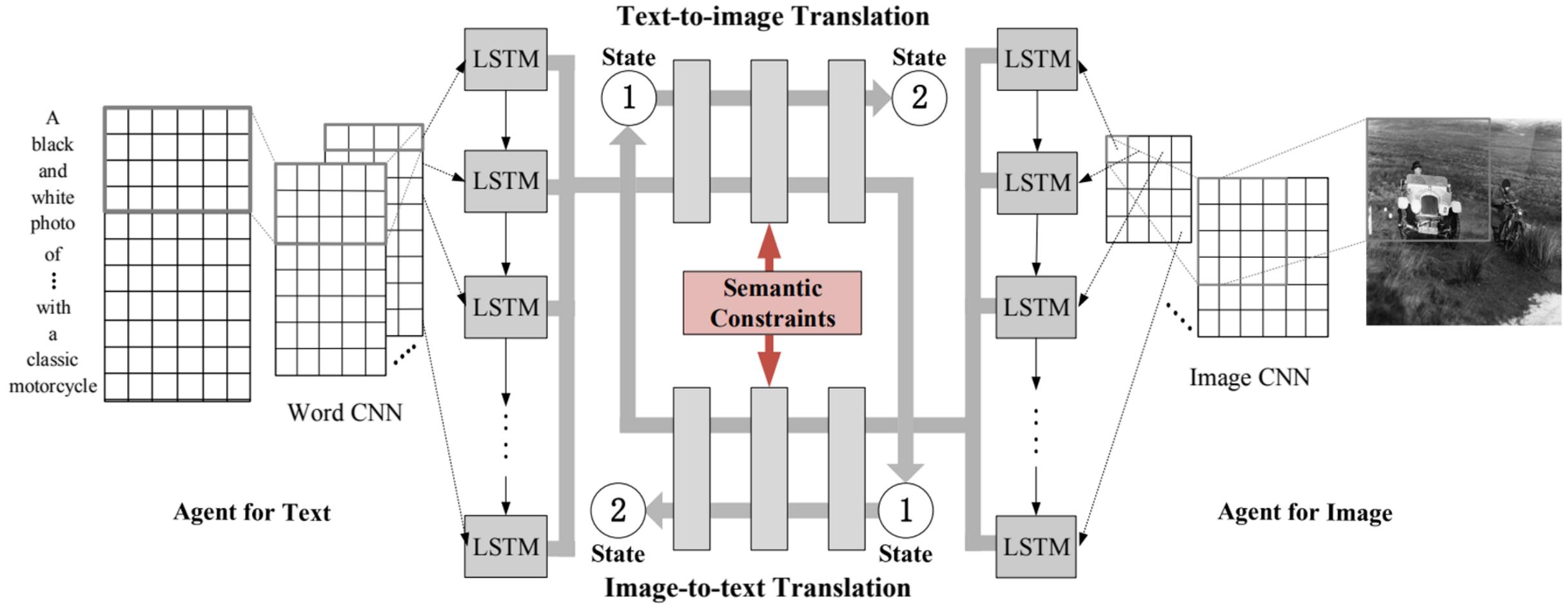


(c) XMediaNet dataset

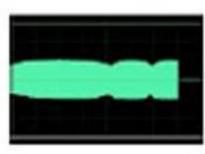
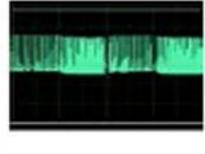
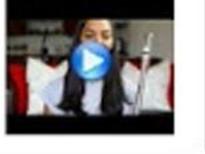
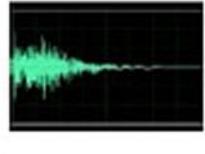
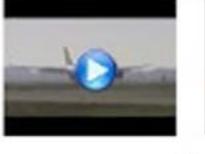
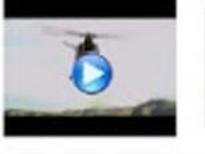
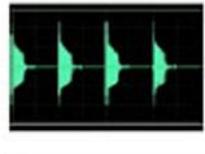
$$r_p = \alpha r_p^{inter} + (1 - \alpha) r_p^{intra}$$

# MODEL

---



# RESULTS

	Image	Text	Audio	Video	3D
violin					
flute					
airplane					
helicopter					
camera					

Method	MAP scores		
	Image→Text	Text→Image	Average
<b>Our CBT Approach</b>	<b>0.577</b>	<b>0.575</b>	<b>0.576</b>
CCL	0.537	0.528	0.533
ACMR	0.536	0.519	0.528
CMDN	0.485	0.516	0.501
Deep-SM	0.399	0.342	0.371
LGCFL	0.441	0.509	0.475
JRL	0.488	0.405	0.447
DCCA	0.425	0.433	0.429
Corr-AE	0.469	0.507	0.488
KCCA	0.252	0.270	0.261
CFA	0.252	0.400	0.326
CCA	0.212	0.217	0.215

▲ PKU XMediaNet Dataset



- 꿀!