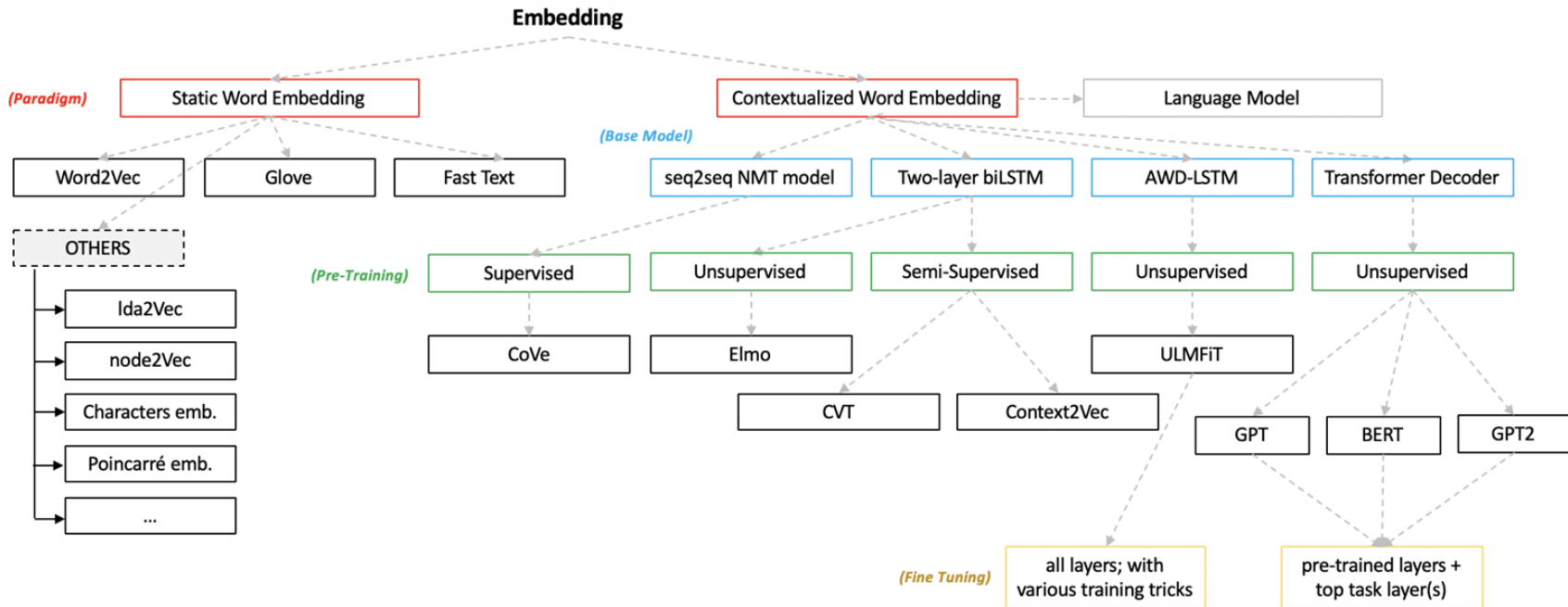
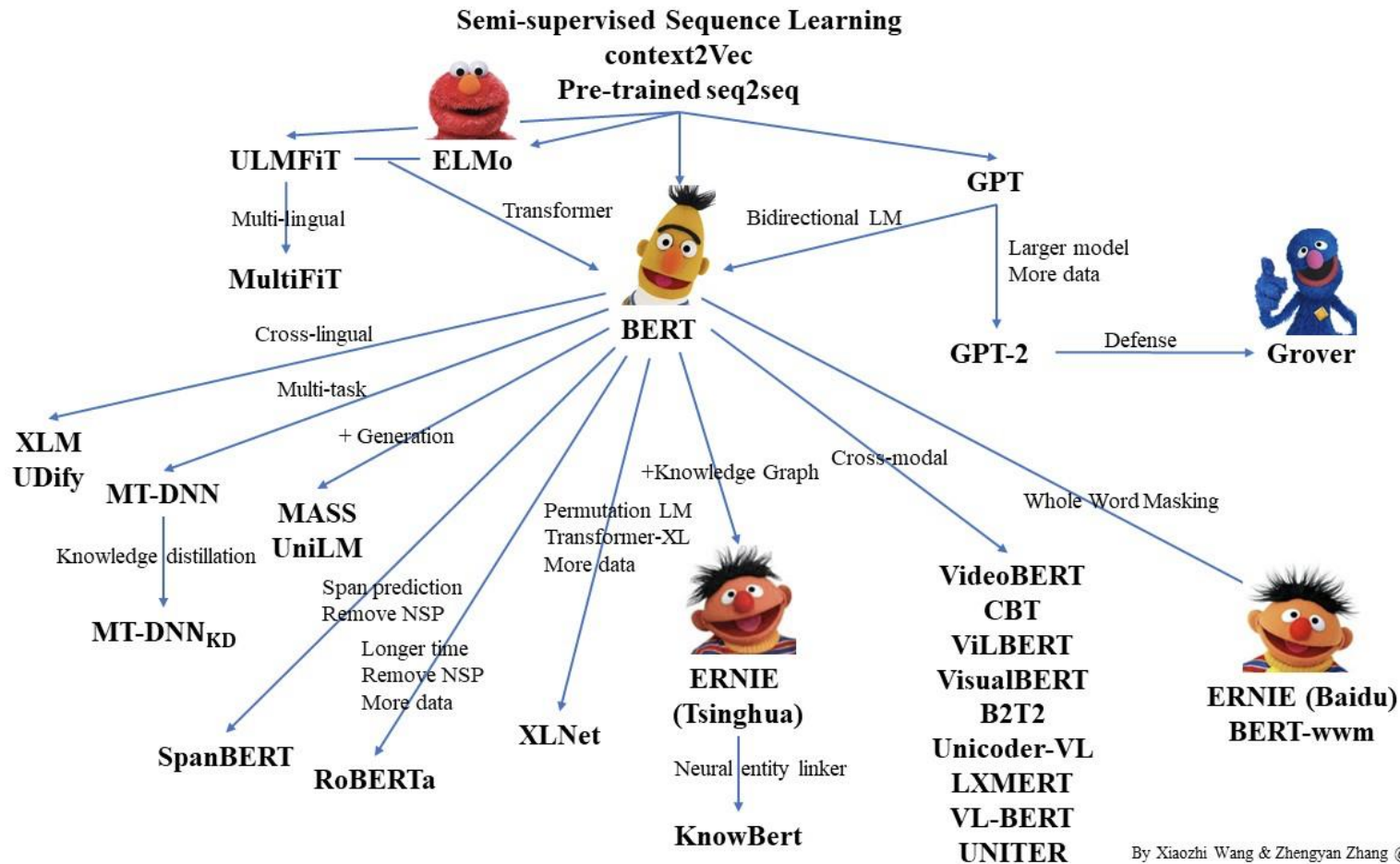


# 문맥적 단어 임베딩 (Contextualized Embedding)

자연어처리 텍스트마이닝





By Xiaozhi Wang & Zhengyan Zhang @THUNLP

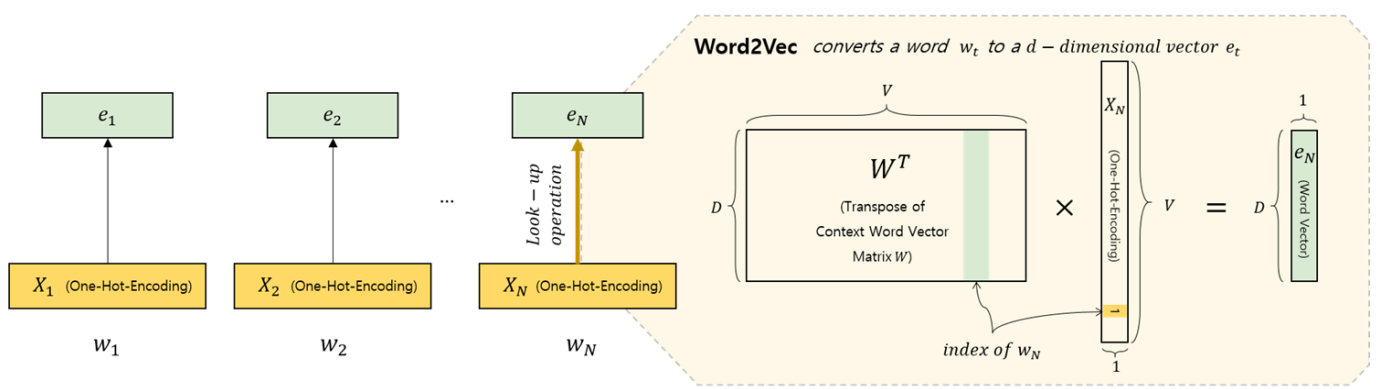
# 문맥적 단어 임베딩 (Contextualized Word Embedding)

---

- Contextualized Word Embedding은 문맥에 따라 vector를 생성
- 같은 단어여도 문맥에 따라 다른 vector가 생성될 수 있음
- 대표적으로 ELMo, BERT, OpenAI GPT
- 이들의 특징은 같은 단어라도 문맥에 따라 다른 방식으로 표현(representation)

# 문맥적 단어 임베딩 (Contextualized Word Embedding)

- Word2Vec의 Embedding은 단어 단위
- 각 단어의 one-hot-encoding vector가  $W_T$ 와 곱해져서 Word Vector를 얻게 됨

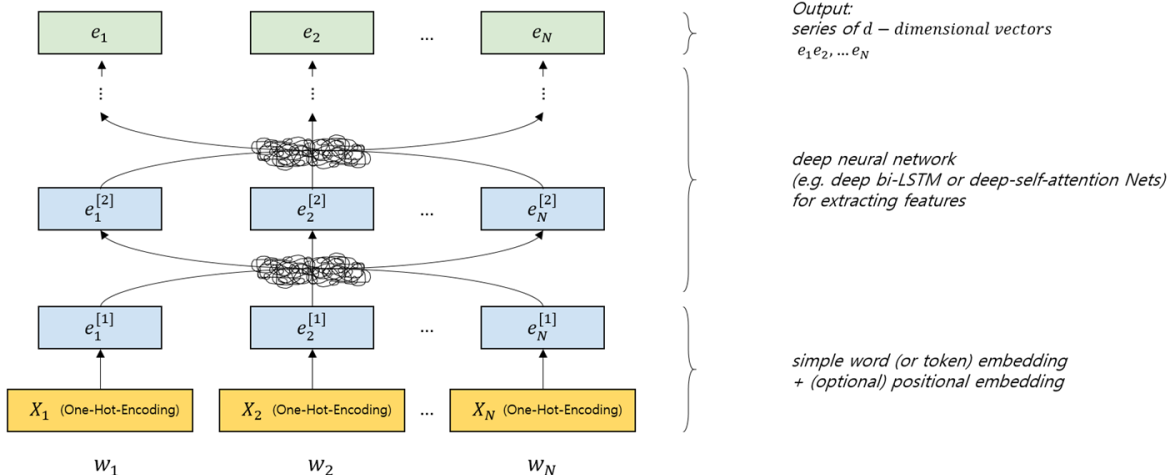


# 문맥적 단어 임베딩 (Contextualized Word Embedding)

- Contextualized Word Representation의 경우 문장을 입력 받아 각 단어에 대한 representation을 산출
- ‘문맥’에 의존적인 ‘단어의 의미’를 잡아내는 feature가 산출

## Contextualized Word Representations

converts a sequence of words (more precisely, tokens)  $w_1 w_2 \dots w_N$  to a series of  $d$ -dimensional vectors  $e_1 e_2, \dots e_N$



# 문맥적 단어 임베딩 (Contextualized Word Embedding)


- Contextualized Word Representation의 경우 문장을 입력 받아 각 단어에 대한 representation을 산출
- ‘문맥’에 의존적인 ‘단어의 의미’를 잡아내는 feature가 산출

관점	Word2Vec	Contextualized Word Representation
Input	단어 단위	문장 단위 (단어의 시퀀스)
Layer	(일반적으로) 단층	(일반적으로) 다계층
Output	해당 단어에 대한 Embedding	문장을 구성하는 각 단어에 대한 Embedding들

# 문맥적 단어 임베딩 (Contextualized Word Embedding)


- 단순한 단어 임베딩 예시

I play piano.

The diagram shows three separate colored squares representing word embeddings for the words 'I', 'play', and 'piano'. The first square is orange, the second is yellow, and the third is blue.

- 순환 신경망 언어 모델을 활용한 단어 임베딩 예시

I → play → piano

The diagram shows three sequential word embeddings for the words 'I', 'play', and 'piano'. The first square is orange. The second square is composed of two overlapping colors, orange and yellow. The third square is composed of three overlapping colors, orange, yellow, and blue.



# 문맥적 단어 임베딩 (Contextualized Word Embedding)

- 좌우 문맥 고려

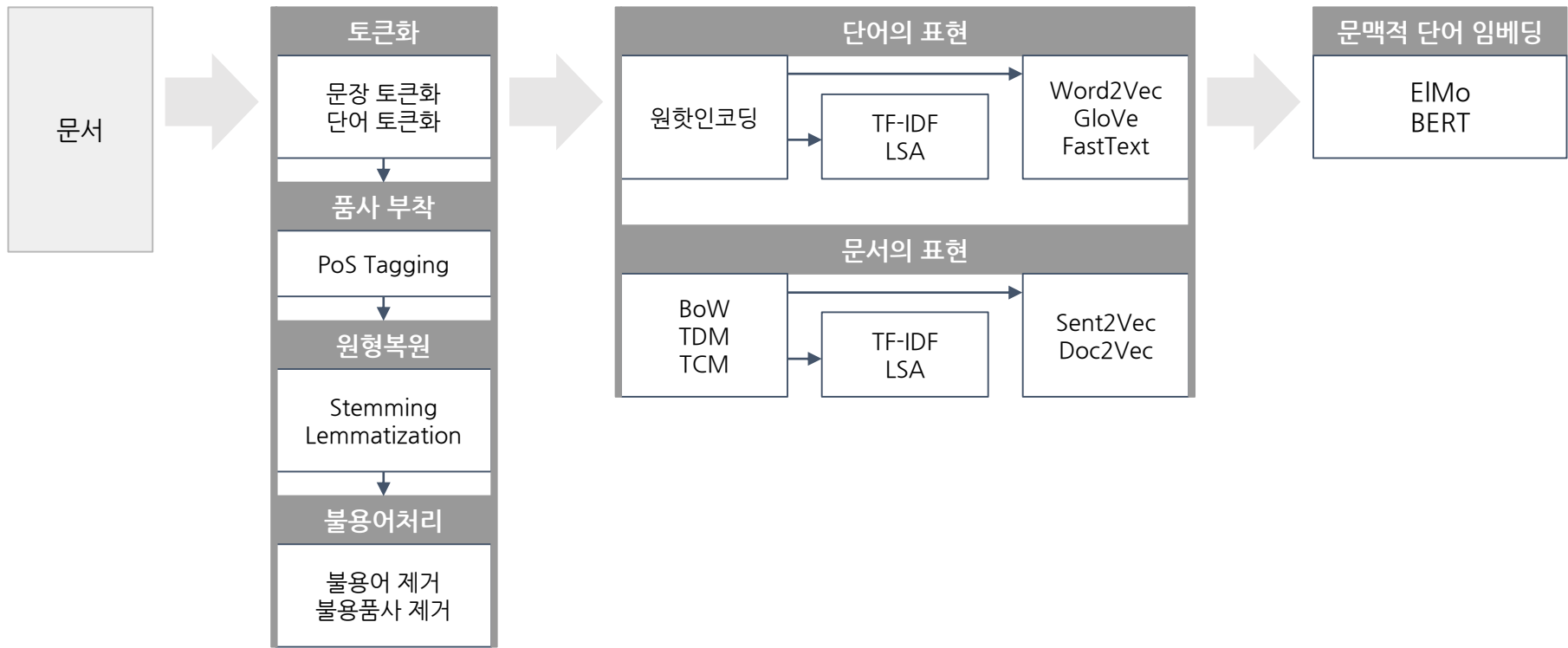
좌 → 우  
I play piano.

I  
I → play  
I → play → piano

우 → 좌  
I play piano.

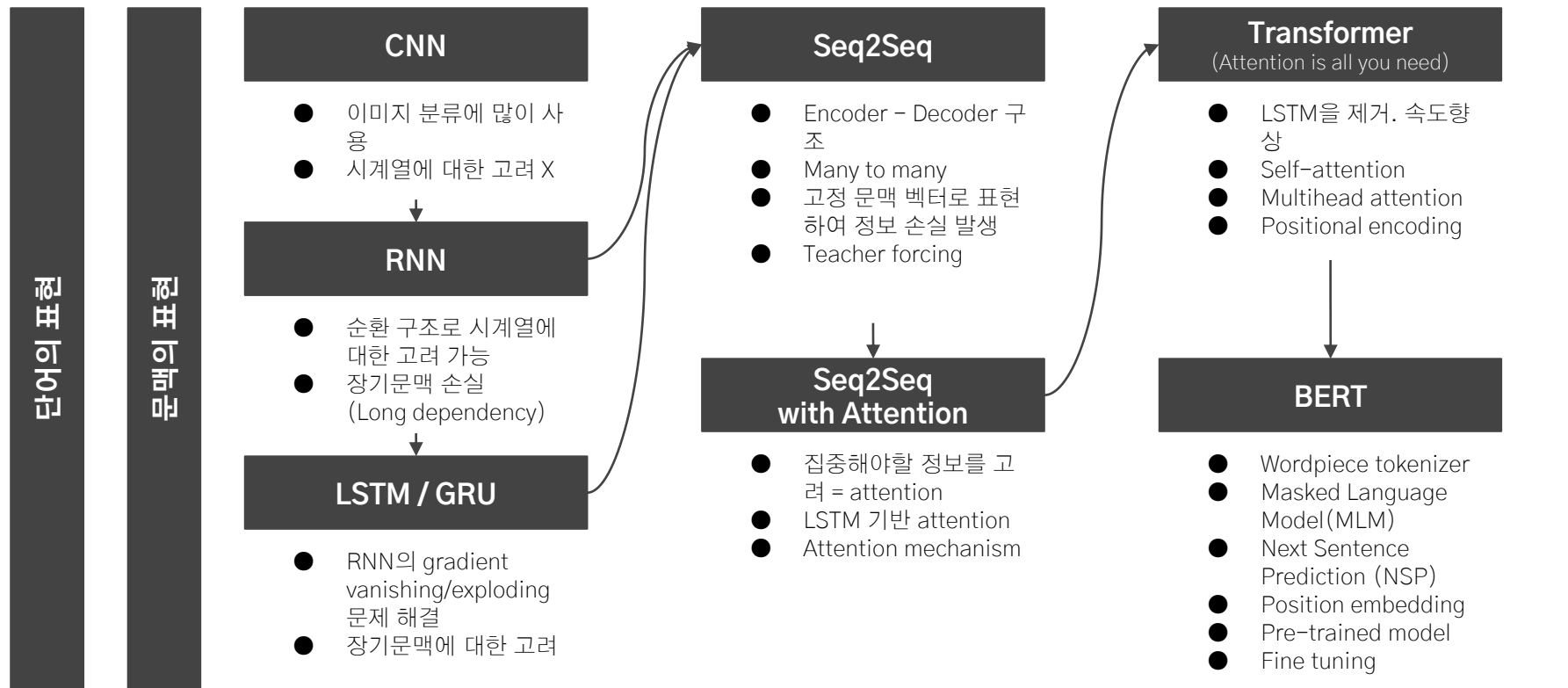
piano  
piano → play  
piano → play → I

# 임베딩 절차



처리 의존도를 고려하여 성능을 높이는 작업을 진행할 수 있음

# BERT 까지



# Tensor

## Tensor Ranks, Shapes, and Types

```
t = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

Rank	Math entity	Python example
0	Scalar (magnitude only)	<code>s = 483</code>
1	Vector (magnitude and direction)	<code>v = [1.1, 2.2, 3.3]</code>
2	Matrix (table of numbers)	<code>m = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]</code>
3	3-Tensor (cube of numbers)	<code>t = [[[2], [4], [6]], [[8], [10], [12]], [[14], [16], [18]]]</code>
n	n-Tensor (you get the idea)	<code>....</code>

# Tensor in NLP

---

sentence	vector representation
hi John	[ [1,0,0,0], [0,1,0,0] ]
hi James	[ [1,0,0,0], [0,0,1,0] ]
hi Brian	[ [1,0,0,0], [0,0,0,1] ]

**(3, 2, 4) 3d tensor!**

hi      John      hi      James      hi      Brian  
 [ [ [1,0,0,0], [0,1,0,0] ], [ [1,0,0,0], [0,0,1,0] ], [ [1,0,0,0], [0,0,0,1] ] ]

# Tensor in Image

