

Paper Seminar

Syntax and Semantics in Language Model Representation

Myeongsup Kim

Integrated M.S./Ph.D. Student
Data Science & Business Analytics Lab.
School of Industrial Management Engineering
Korea University

Myeongsup_kim@korea.ac.kr

A Structural Probe for Finding Syntax in Word Representation

Hewitt and Manning, 2019, NAACL

Visualizing and Measuring the Geometry of BERT

Coenen et al., 2019, NIPS

Introduction

- **Concept of Language Model**

Introduction

-What This Seminar Does Not Cover

<What This Seminar Does Not Cover>

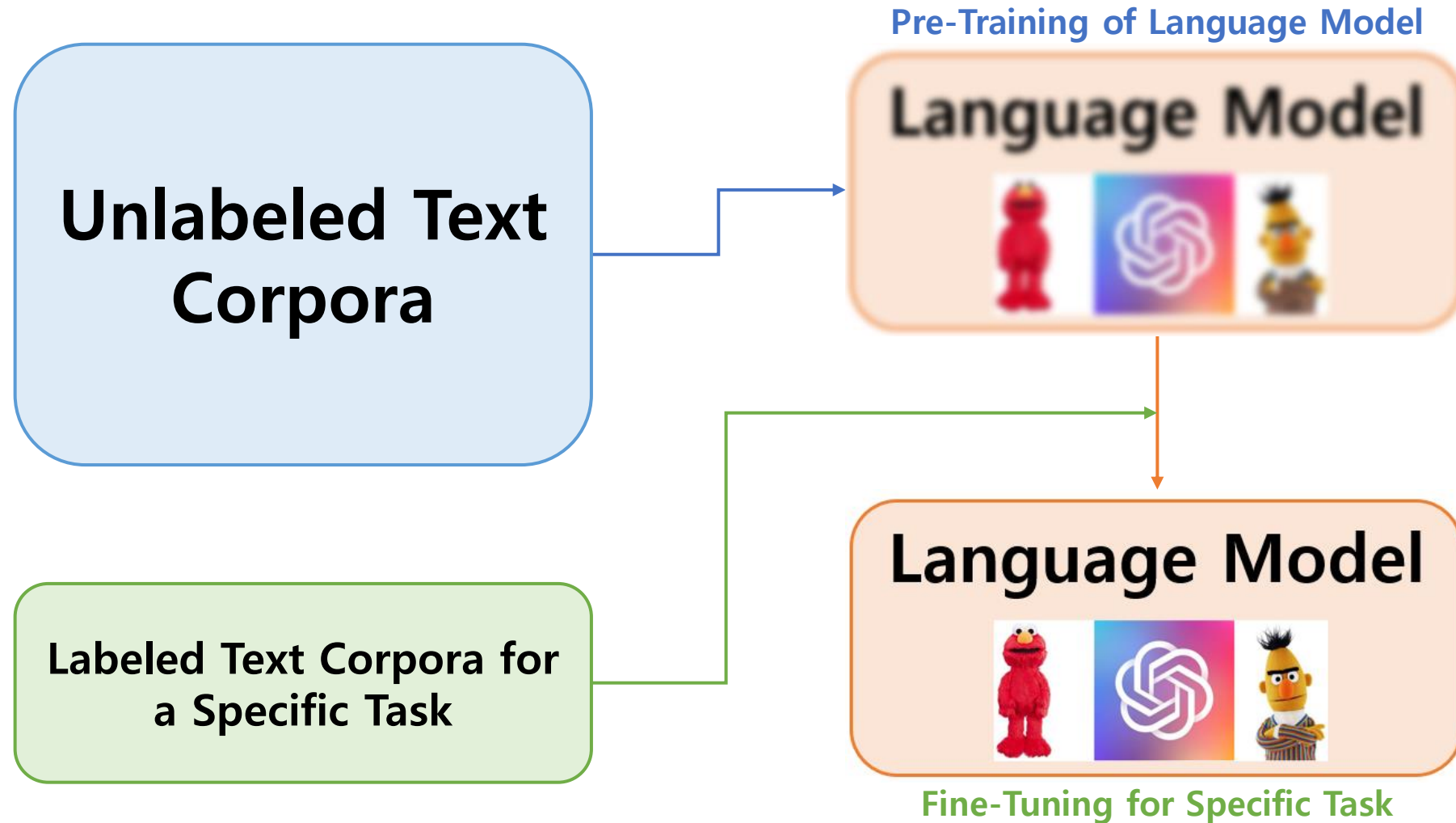
- **Details of ELMo**

[Peters et al., 2018, Deep Contextualized Word Representations, NAACL](#)

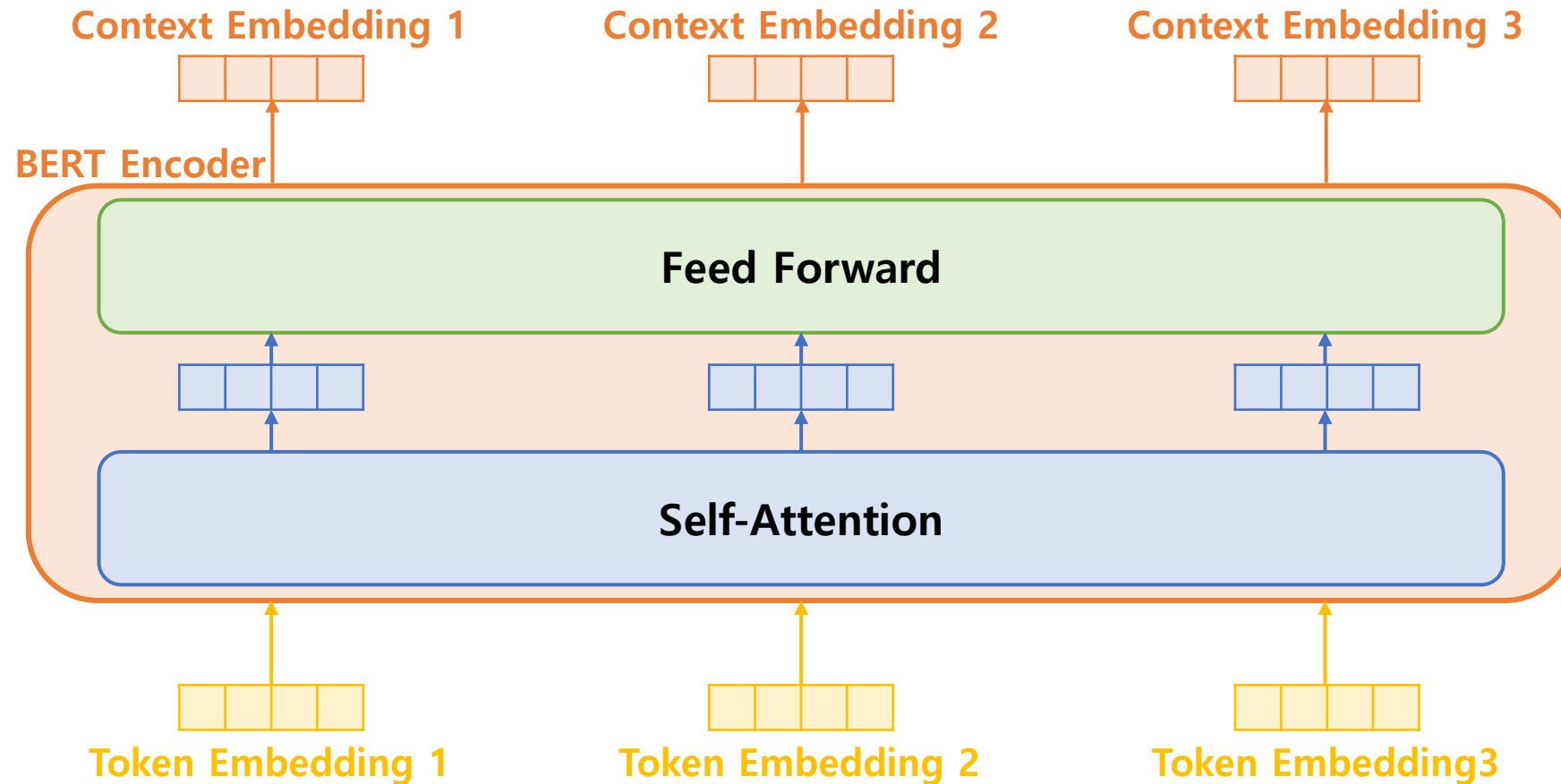
- **Details of BERT**

[Devlin et al., 2019, BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding, NAACL](#)

<Concept of Language Model : Motivation>



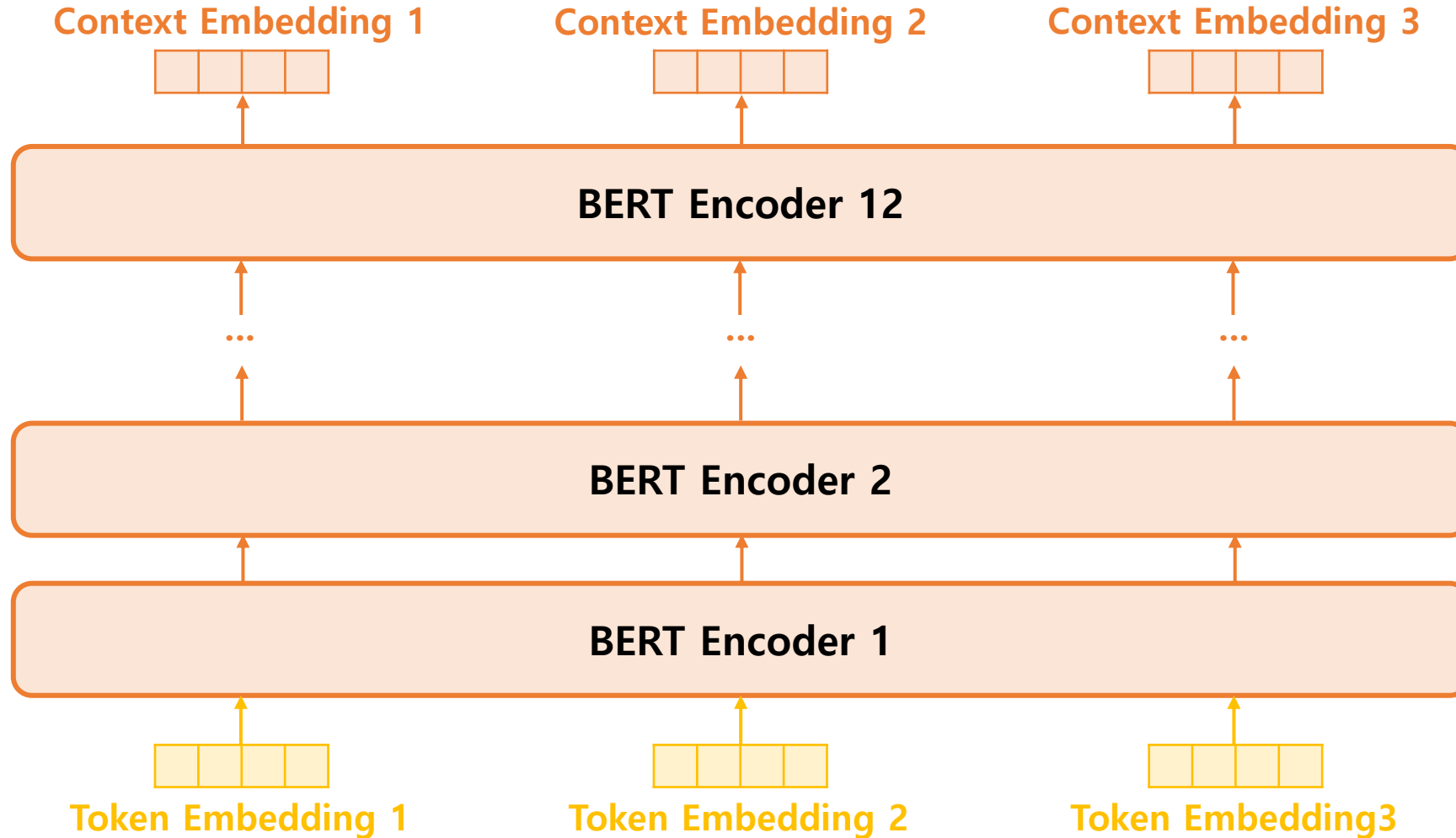
<Concept of Language Model : Pre-Training>



Introduction

-Concept of Language Model

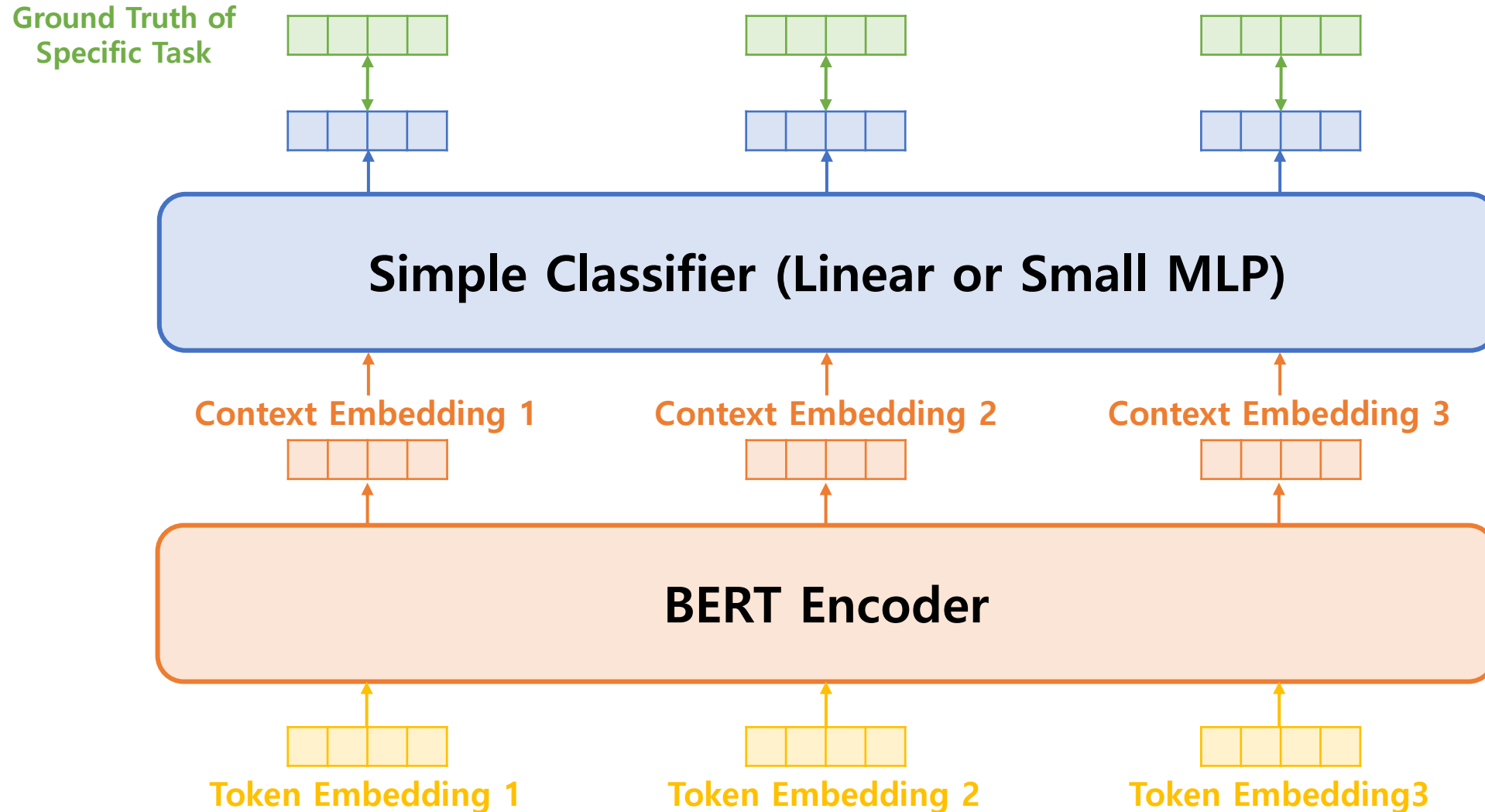
<Concept of Language Model : Pre-Training>



Introduction

-Concept of Language Model

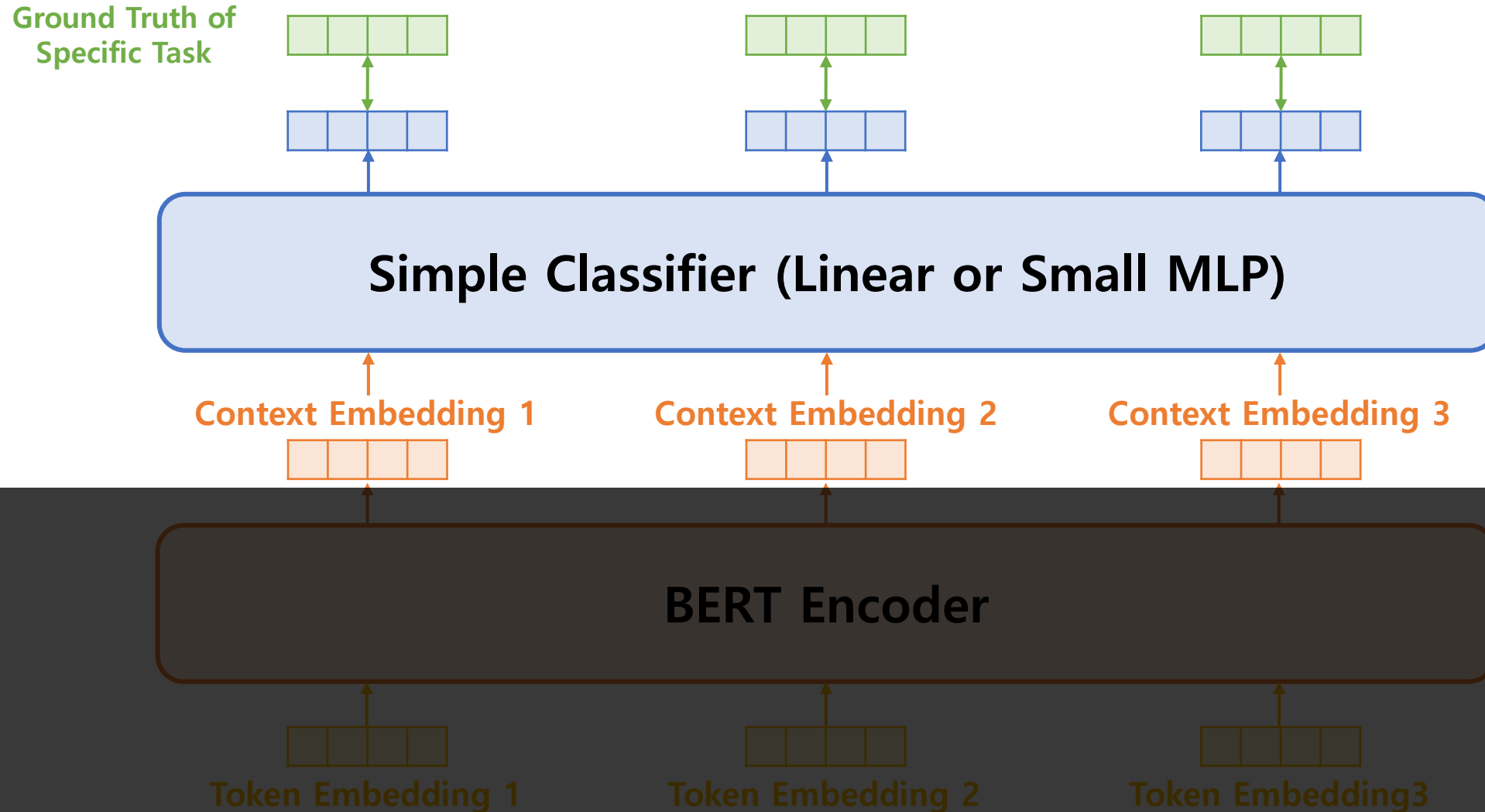
<Concept of Language Model : Fine-tuning>



Introduction

-Concept of Language Model

<Concept of Language Model : Fine-tuning>



Introduction

-Concept of Language Model

<Concept of Language Model : Fine-tuning>

Ground Truth of
Specific Task

So, What Does Language Model Learn from the Context of
Unlabeled Text Corpora?

Which Information is Contained in the Context Embeddings?

Context Embedding 1

Context Embedding 2

Context Embedding 3

BERT Encoder

Token Embedding 1

Token Embedding 2

Token Embedding 3

A Structural Probe for Finding Syntax in Word Representation

Hewitt and Manning, 2019, NAACL

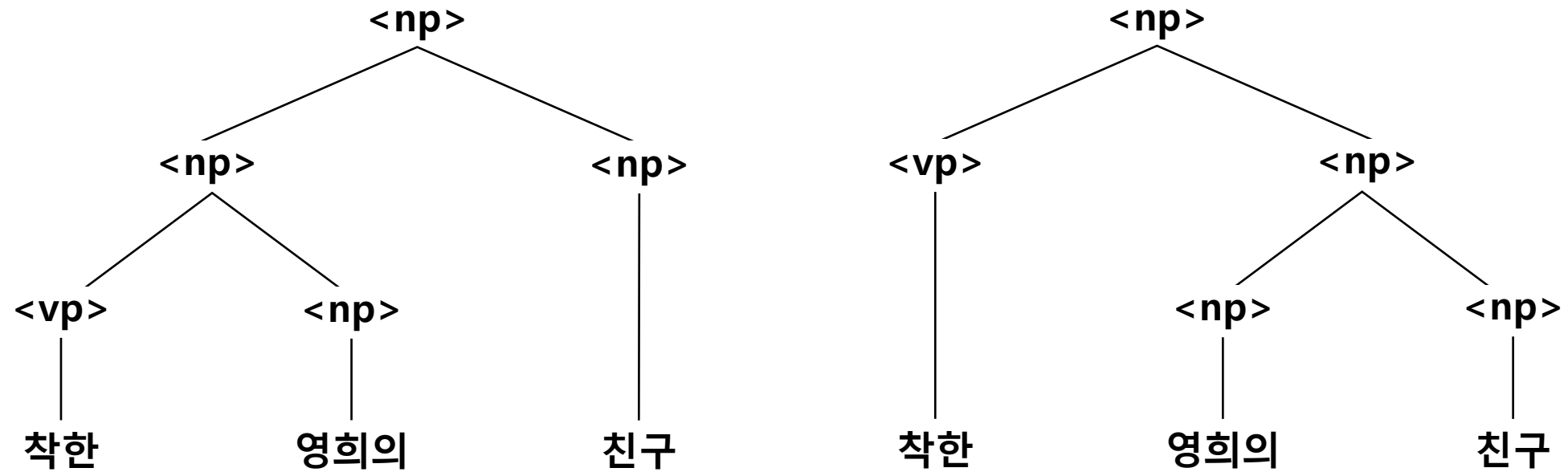
A Structural Probe for Finding Syntax in Word Representation

Hewitt and Manning, 2019, NAACL

Pre-requisites

- **Dependency Parse Tree**

<Dependency Parse Tree>



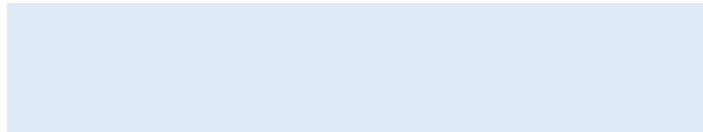
Pre-requisites

-Dependency Parse Tree

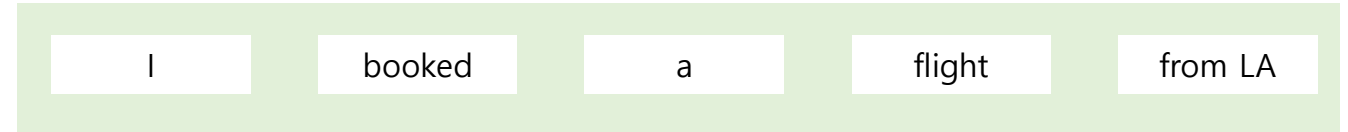
<Transition-based Dependency Parsing>

- Buffer의 제일 앞에 있는 Word 하나를 Stack으로 옮긴 후(Shift)
- Stack에 있는 단어에 대해서 문법 규칙을 적용하고 하나의 Word만 남기는 행동(Reduce)을
- Stack에 하나의 Word만 남을 때 까지 반복하는 방식으로 Parsing
 - Buffer : 토큰나이징 한 단어들 중에 Parsing을 기다리고 있는 단어들
 - Stack : 토큰나이징 한 단어들 중에 파싱 작업을 진행 중인 단어들

Stack



Buffer



root

I

booked

a

flight

from LA

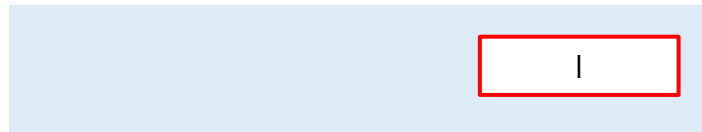
Pre-requisites

-Dependency Parse Tree

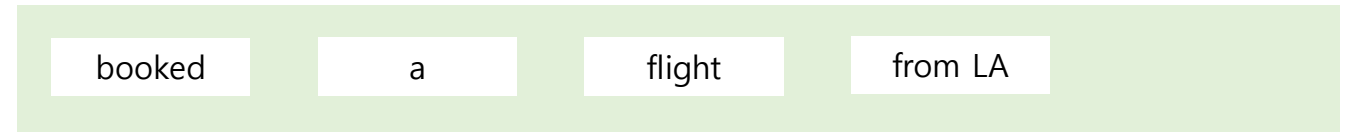
<Transition-based Dependency Parsing>

- Buffer의 제일 앞에 있는 Word 하나를 Stack으로 옮긴 후(Shift)
- Stack에 있는 단어에 대해서 문법 규칙을 적용하고 하나의 Word만 남기는 행동(Reduce)을
- Stack에 하나의 Word만 남을 때 까지 반복하는 방식으로 Parsing
 - Buffer : 토큰나이징 한 단어들 중에 Parsing을 기다리고 있는 단어들
 - Stack : 토큰나이징 한 단어들 중에 파싱 작업을 진행 중인 단어들

Stack



Buffer



Shift

root

I

booked

a

flight

from LA

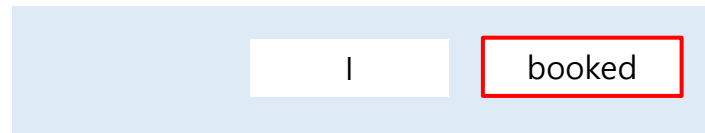
Pre-requisites

-Dependency Parse Tree

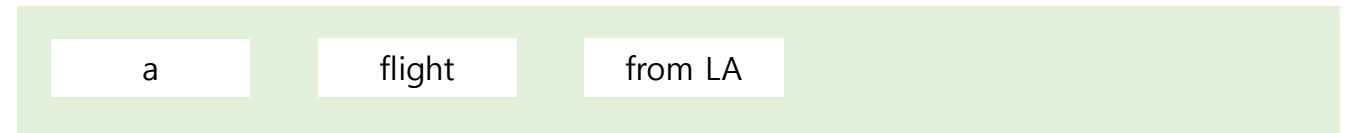
<Transition-based Dependency Parsing>

- Buffer의 제일 앞에 있는 Word 하나를 Stack으로 옮긴 후(Shift)
- Stack에 있는 단어에 대해서 문법 규칙을 적용하고 하나의 Word만 남기는 행동(Reduce)을
- Stack에 하나의 Word만 남을 때 까지 반복하는 방식으로 Parsing
 - Buffer : 토큰나이징 한 단어들 중에 Parsing을 기다리고 있는 단어들
 - Stack : 토큰나이징 한 단어들 중에 파싱 작업을 진행 중인 단어들

Stack



Buffer



Shift

root I booked a flight from LA

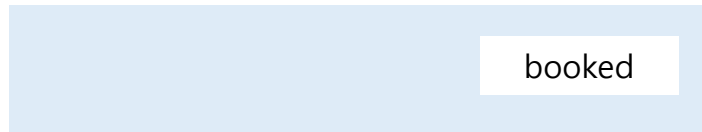
Pre-requisites

-Dependency Parse Tree

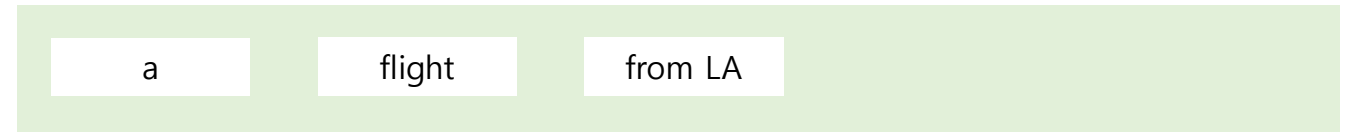
<Transition-based Dependency Parsing>

- Buffer의 제일 앞에 있는 Word 하나를 Stack으로 옮긴 후(Shift)
- Stack에 있는 단어에 대해서 문법 규칙을 적용하고 하나의 Word만 남기는 행동(Reduce)을
- Stack에 하나의 Word만 남을 때 까지 반복하는 방식으로 Parsing
 - Buffer : 토큰나이징 한 단어들 중에 Parsing을 기다리고 있는 단어들
 - Stack : 토큰나이징 한 단어들 중에 파싱 작업을 진행 중인 단어들

Stack



Buffer



Left-Reduce



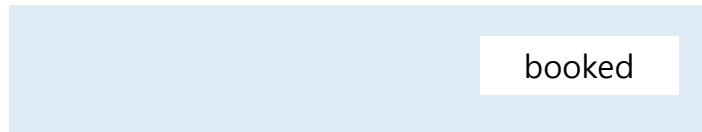
Pre-requisites

-Dependency Parse Tree

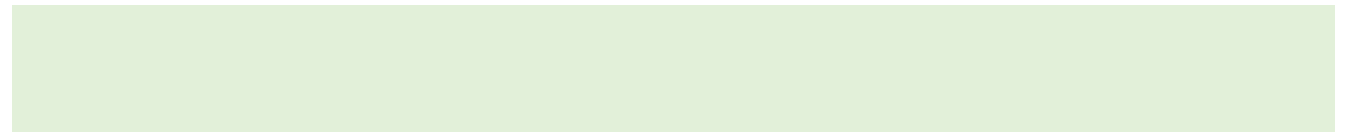
<Transition-based Dependency Parsing>

- Buffer의 제일 앞에 있는 Word 하나를 Stack으로 옮긴 후(Shift)
- Stack에 있는 단어에 대해서 문법 규칙을 적용하고 하나의 Word만 남기는 행동(Reduce)을
- **Stack에 하나의 Word만 남을 때 까지 반복하는 방식으로 Parsing**
 - Buffer : 토큰나이징 한 단어들 중에 Parsing을 기다리고 있는 단어들
 - Stack : 토큰나이징 한 단어들 중에 파싱 작업을 진행 중인 단어들

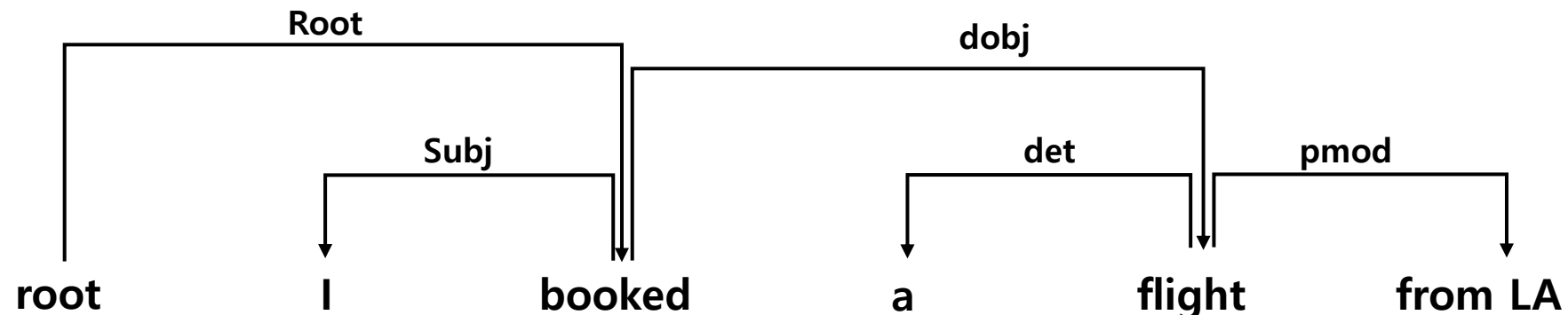
Stack



Buffer



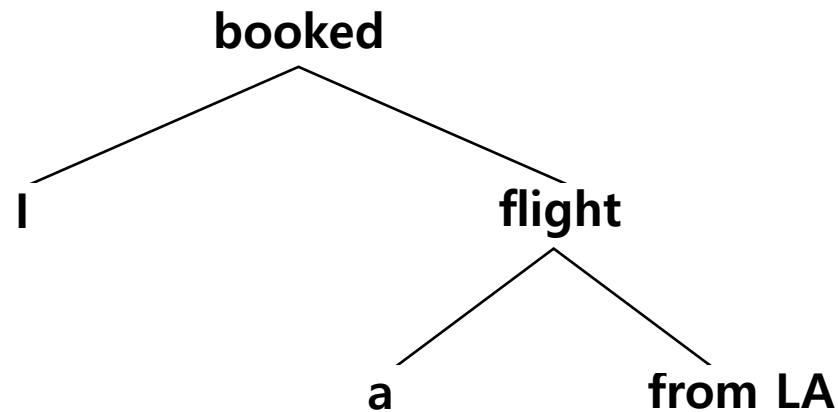
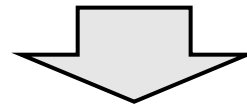
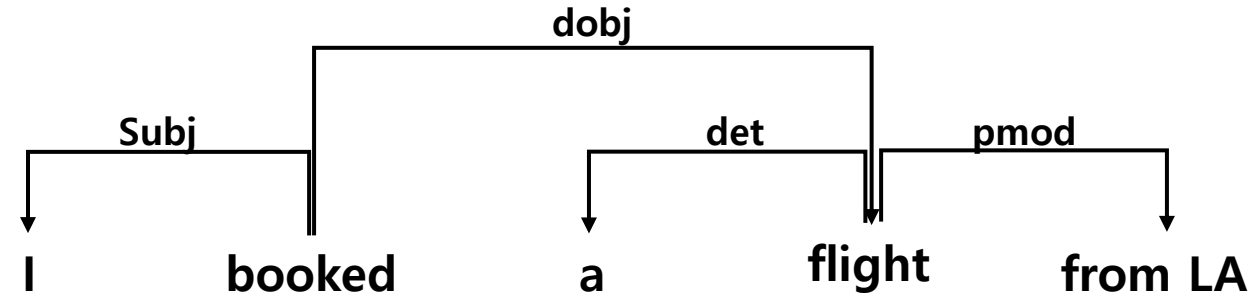
Final State



Pre-requisites

-Dependency Parse Tree

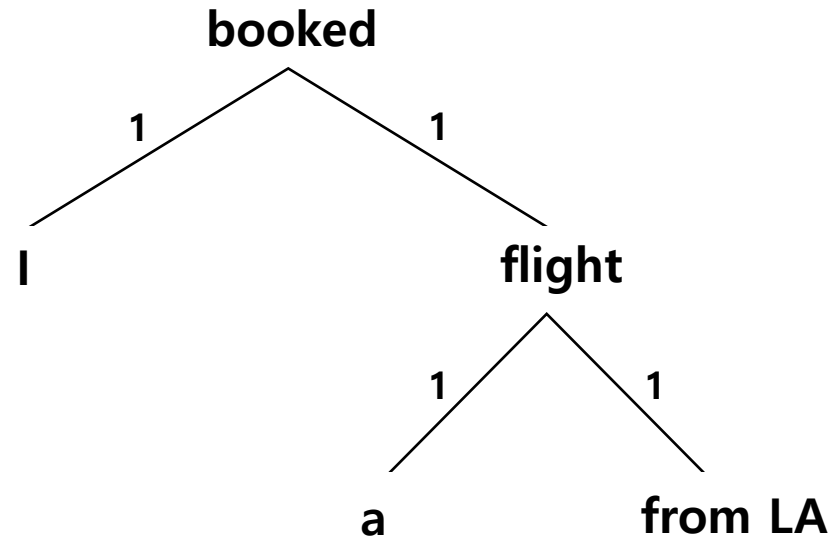
<Dependency Parse Tree>



Pre-requisites

-Dependency Parse Tree

<Dependency Parse Tree>



$$d_T(\text{booked}, I) = 1$$

$$d_T(I, a) = 3$$

$$d_T(\text{booked}, \text{flight}) = 1$$

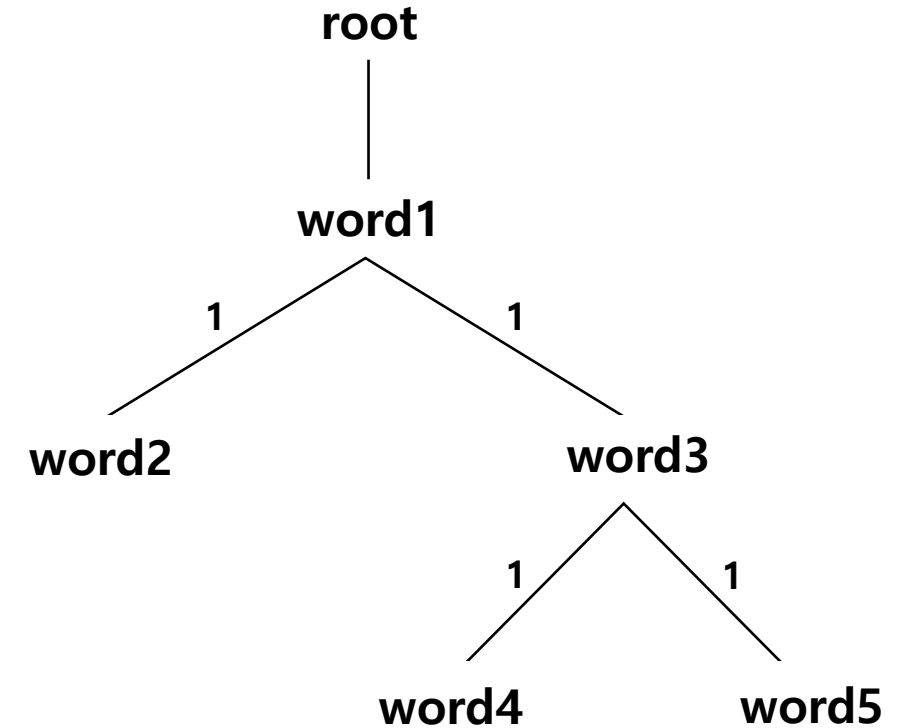
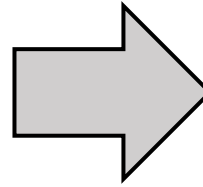
Pre-requisites

-Dependency Parse Tree

<Dependency Parse Tree>

for i, j : $d_T(\text{word}_i, \text{word}_j)$

for all i : $d_T(\text{root}, \text{word}_i)$



<Unique Tree>

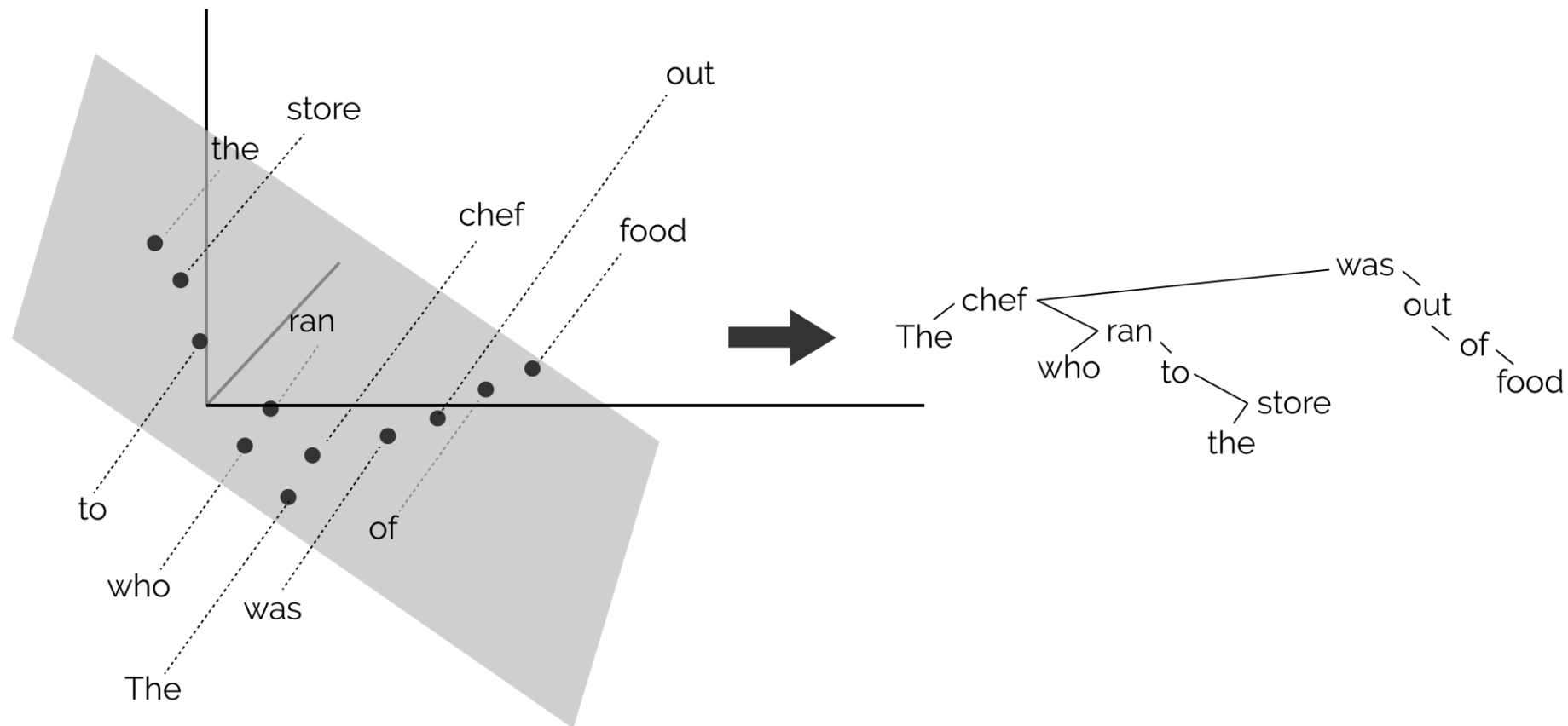
Method

- **Structural Probe**

<Structural Probe Overview>

There is a Linear Transformation that Transforms the Embeddings of Language Model into Dependency Parse Tree

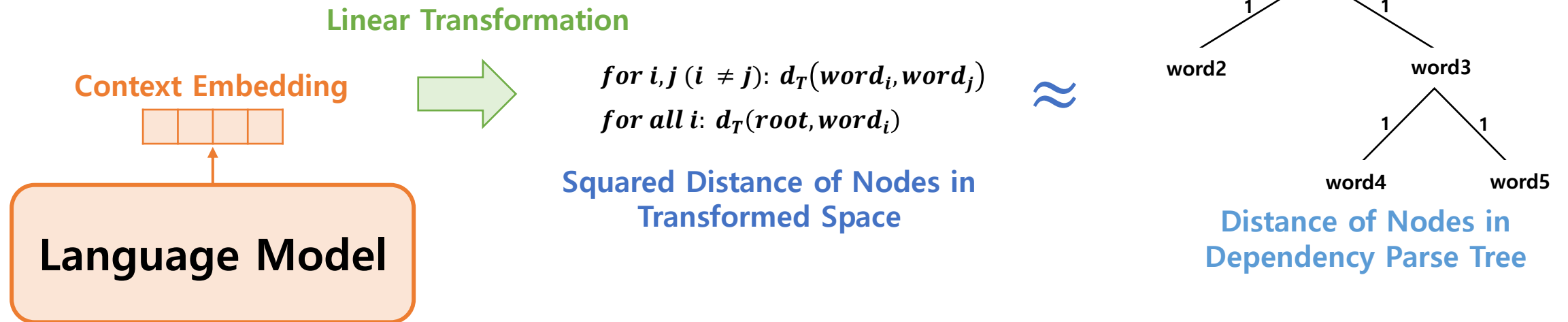
So, Language Model Embeds Dependency Relation (Syntactic Information) of Text



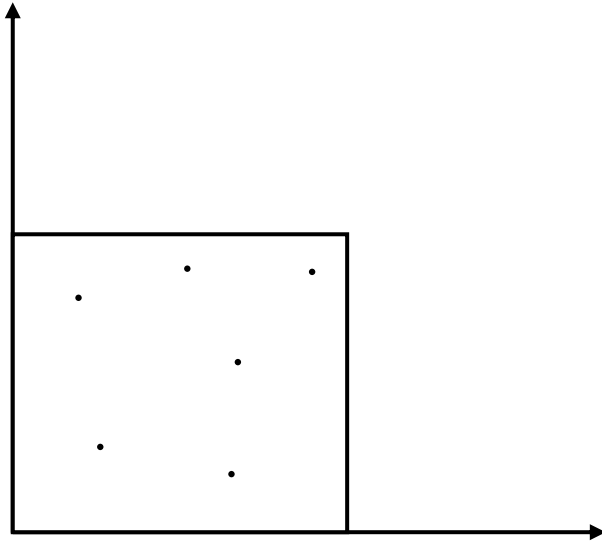
Method

-Structural Probe

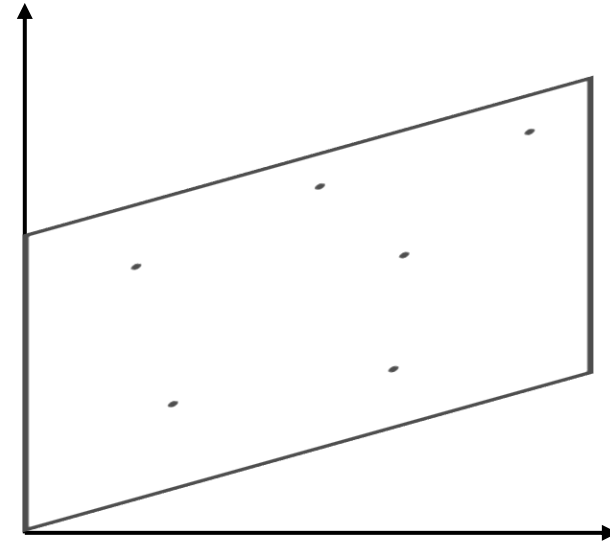
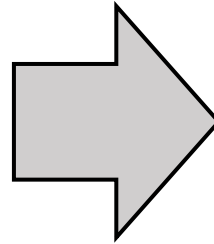
<Structural Probe Overview>



<Geometric Meaning of Linear Transformation>



$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$



$$Av = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

Method

-Structural Probe

<Structural Probe>

<Notation>

$w_{1:n}^l$: words in sequence_l

$h_{1:n}^l$: sequence of vector representation

A : positive semi definite, symmetric matrix, $A \in \mathbb{S}_+^{m \times m}$

$h^T A h$: family of inner product

B : linear transformation, $B \in \mathbb{R}^{k \times m}$, such that $A = B^T B$

d_m : embedding dimension

$$h = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_m \end{bmatrix} (m \times 1) \quad A = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{1k} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2k} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kk} & \cdots & a_{km} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mk} & \cdots & a_{mm} \end{bmatrix} (m \times m) \quad B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1m} \\ b_{21} & b_{22} & \cdots & b_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ b_{k1} & b_{k2} & \cdots & b_{km} \end{bmatrix} (k \times m)$$

$h^T A h$: scalar (1×1)

<Linear Transformation>

$h^T A h$: family of inner product

B : linear transformation, $B \in \mathbb{R}^{k \times m}$, such that $A = B^T B$

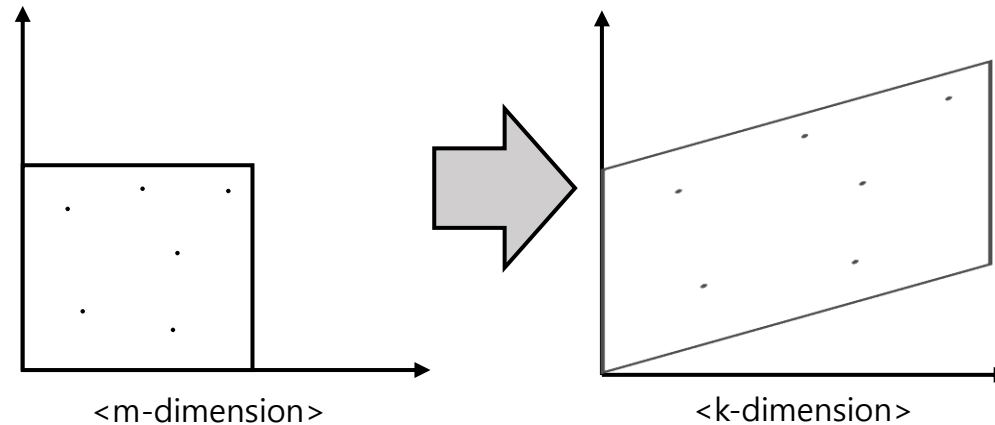
$\therefore h^T A h = (Bh)^T (Bh)$: inner product of transformed vector

$$B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1m} \\ b_{21} & b_{22} & \cdots & b_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ b_{k1} & b_{k2} & \cdots & b_{km} \end{bmatrix} (k \times m)$$

$$h = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_m \end{bmatrix} (m \times 1)$$

nd_k : new dimension

$$Bh = \begin{bmatrix} nd_1 \\ nd_2 \\ \vdots \\ nd_k \end{bmatrix} (k \times 1), k \leq m$$

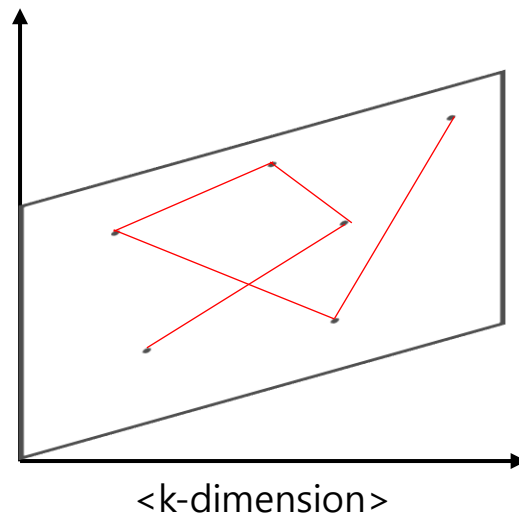


<Calculate Distance of Nodes>

$$d_B(h_i^l, h_j^l)^2 = \left(B(h_i^l - h_j^l) \right)^T (B(h_i^l - h_j^l))$$

$$d_B(h_i^l, h_j^l)^2 = (nd_{i1} - nd_{j1})^2 + (nd_{i2} - nd_{j2})^2 + \dots + (nd_{ik} - nd_{jk})^2$$

$\therefore d_B(h_i^l, h_j^l)^2$: *squared euclidean distance of nodes*



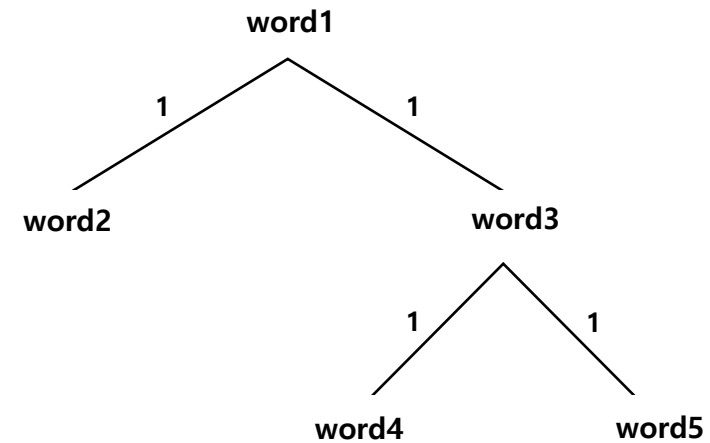
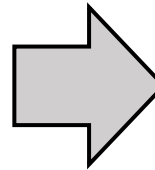
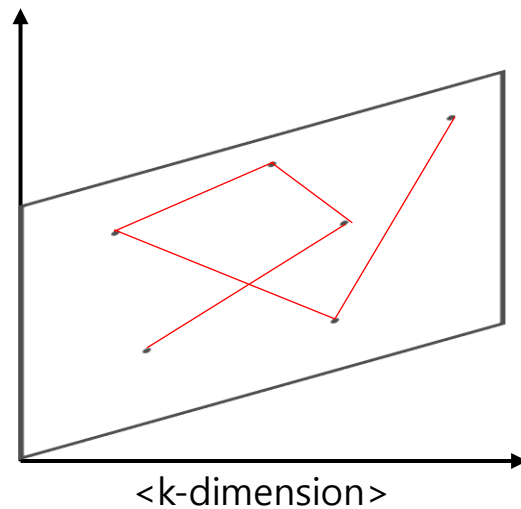
Method

-Structural Probe

<Training>

$$\min_B \sum_l \frac{1}{|s^l|^2} \sum_{i,j} |d_{T^l}(w_i^l, w_j^l) - d_B(h_i^l, d_j^l)|^2$$

$|s^l|$: length of the sentence, sentence has $|s^l|^2$ word pairs

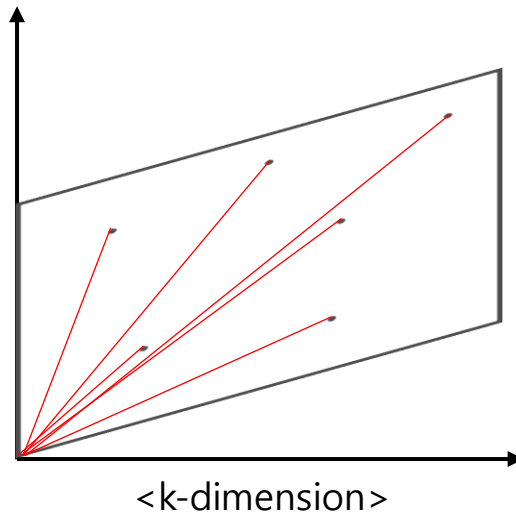


<Calculate Distance from Root Node>

$$||h_i||_B^2 = (Bh_i)^T (Bh_i)$$

$$||h_i||_B^2 = (nd_{i1})^2 + (nd_{i2})^2 + \dots + (nd_{ik})^2$$

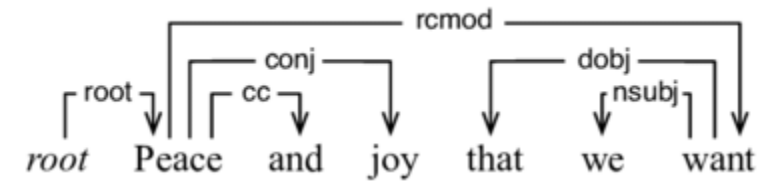
$\therefore ||h_i||_B^2$: *norm of node i in transformed space*



Experiments

<Penn Treebank Dataset>

"Peace and joy that we want"



<Sentence>

<Tree Structure of Given Sentence>

<Dependency Relations of Tokens in Given Sentence>

Experiments

-Representation Models

<Representation Models>



<ELMo>

5.5B-word
Pre-trained ELMo
1024-dim Embedding



<BERTBase>

Pre-trained BERTBase
(cased)
768-dim Embedding



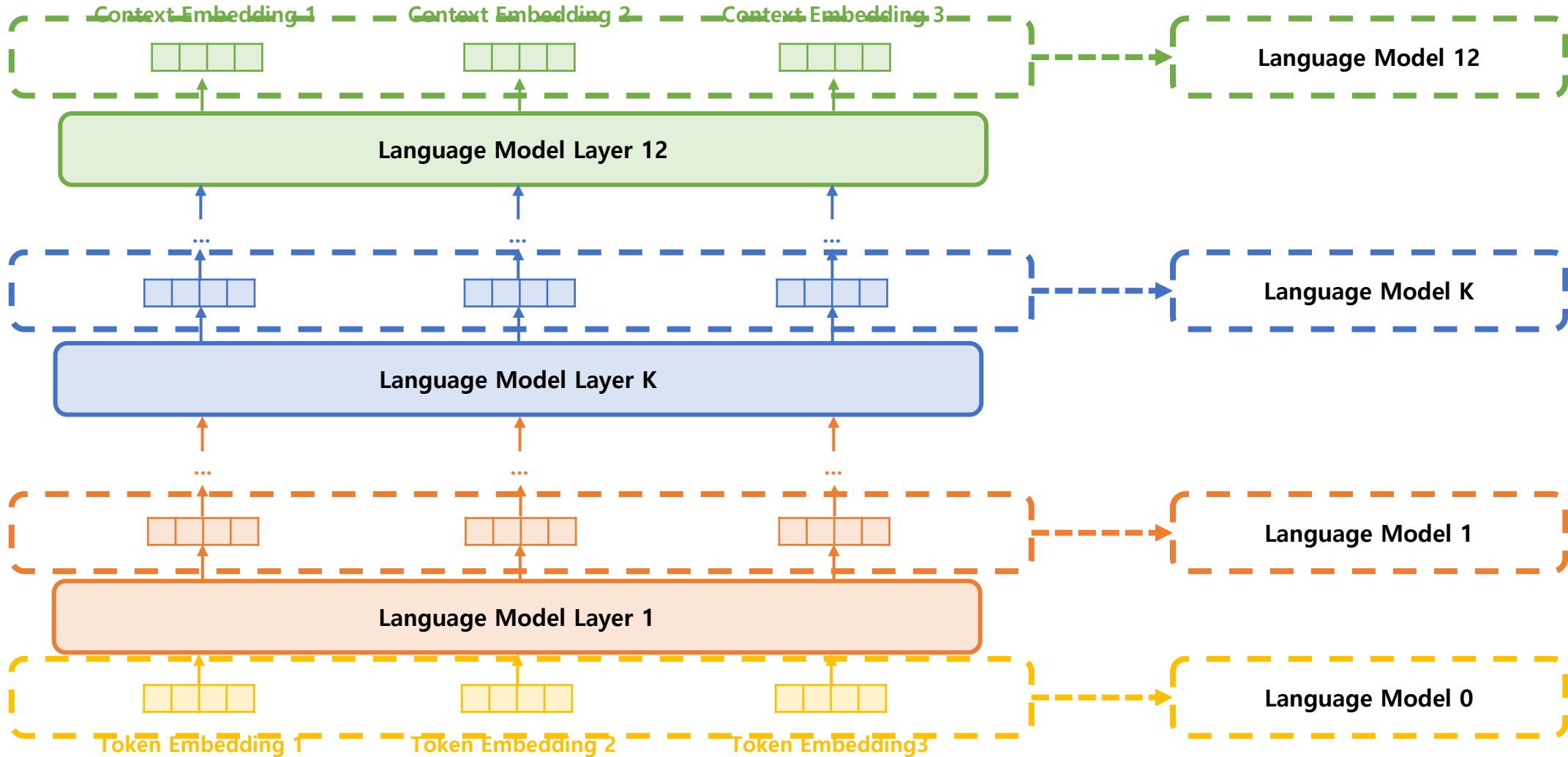
<BERTLarge>

Pre-trained BERTLarge
(cased)
1024-dim Embedding

Experiments

-Representation Models

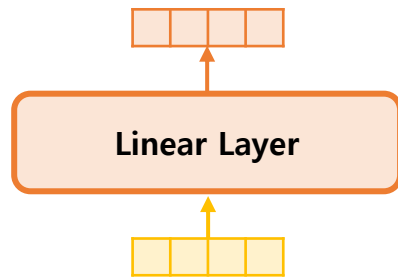
<Representation Models>



Experiments

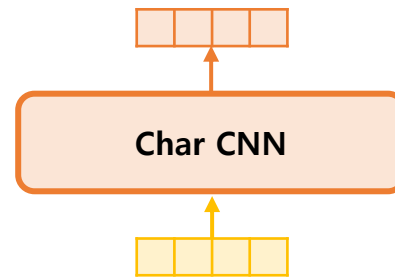
-Baseline

<Baseline>



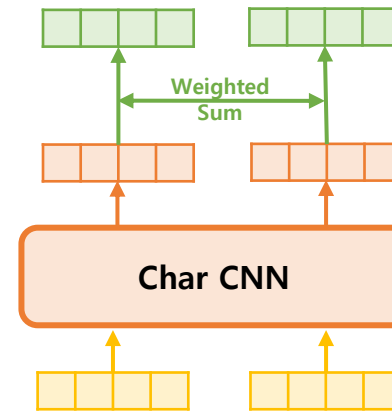
<Linear>

Token Embedding
No Context Information
No Sequential Information



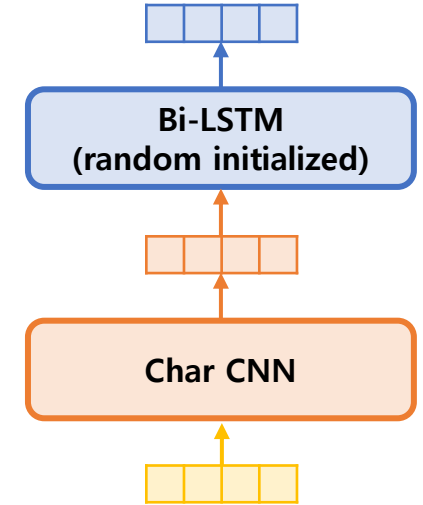
<ELMo0>

Character Embedding
No Context Information
No Sequential Information



<Decay0>

ELMo0 with Weighted Sum of
Other Tokens in Same Sentence
 $Weight = \frac{1}{2^d}, d: \text{distance of words}$
Contain Context Information
No Sequential Information



<Proj0>

ELMo0 with Randomly Initialized
Bi-LSTM (1024-dim)
Contain Sequential and Context
Information

Experiments

-Evaluation Metric

<Evaluation Metric>

<UUAS>

$$\frac{\# \text{ of Correct Nodes}}{\# \text{ of All Groud Truth Nodes}}$$

Root: Node with Least Depth

<Root%>

$$\frac{\# \text{ of Correct Root}}{\# \text{ of Root Node of Ground Truth}}$$

<Distance Metric>

<DSpr.>

$$\frac{\sum_{i=5}^{50} \frac{\text{avg}(\text{spearman Correlation of Each (Predicted, True) word in lenght } i \text{ sentence})}{\# \text{ of length } i \text{ sentence}}}{\sum_{i=5}^{50} i}$$

Ordered by Distance

<Depth Metric>

<NSpr.>

$$\frac{\sum_{i=5}^{50} \frac{\text{avg}(\text{spearman Correlation of Each (Predicted, True) word in lenght } i \text{ sentence})}{\# \text{ of length } i \text{ sentence}}}{\sum_{i=5}^{50} i}$$

Ordered by Norm

Experiments

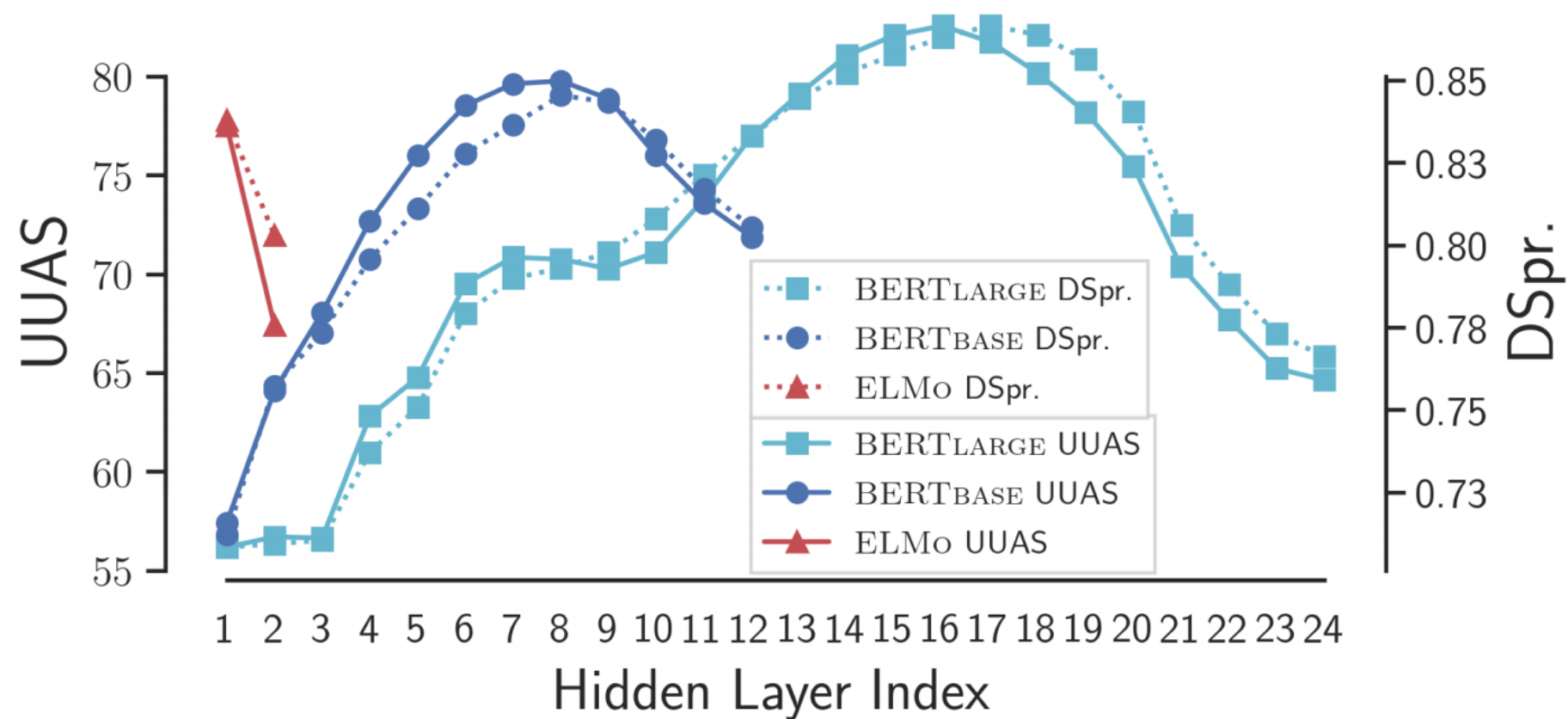
-Result

<Result>

Method	Distance		Depth	
	UUAS	DSpr.	Root%	NSpr.
LINEAR	48.9	0.58	2.9	0.27
ELMO0	26.8	0.44	54.3	0.56
DECAY0	51.7	0.61	54.3	0.56
PROJ0	59.8	0.73	64.4	0.75
ELMO1	77.0	0.83	86.5	0.87
BERT ^{BASE} 7	79.8	0.85	88.0	0.87
BERT ^{LARGE} 15	82.5	0.86	89.4	0.87
BERT ^{LARGE} 16	81.7	0.87	90.1	0.89

<Result of Structural Probe on PTB>

<Result>

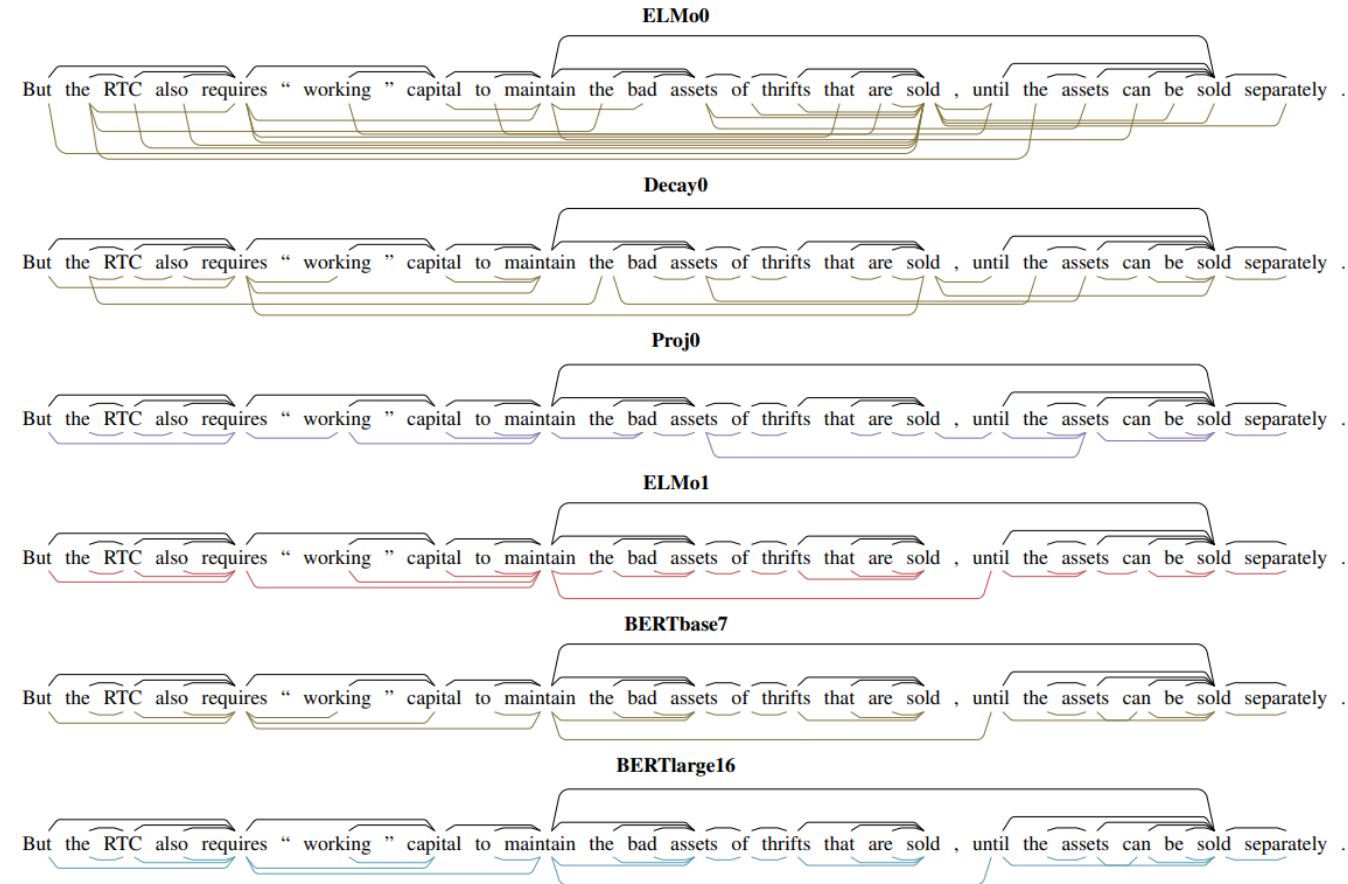


<Parse Distance UUAS & Dspr. Across BERT and ELMo Layers>

Experiments

-Result

<Result>

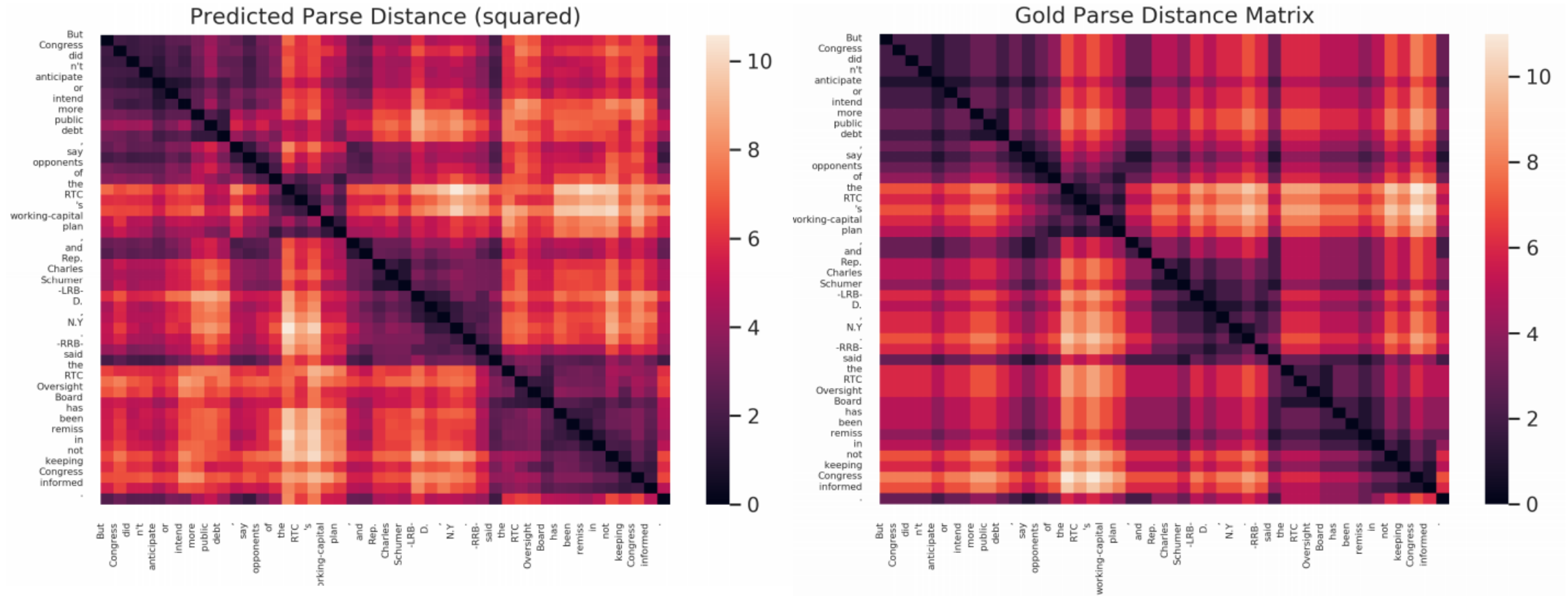


<Minimum Spanning Trees Extracted by Various Models>

Experiments

-Result

<Result>

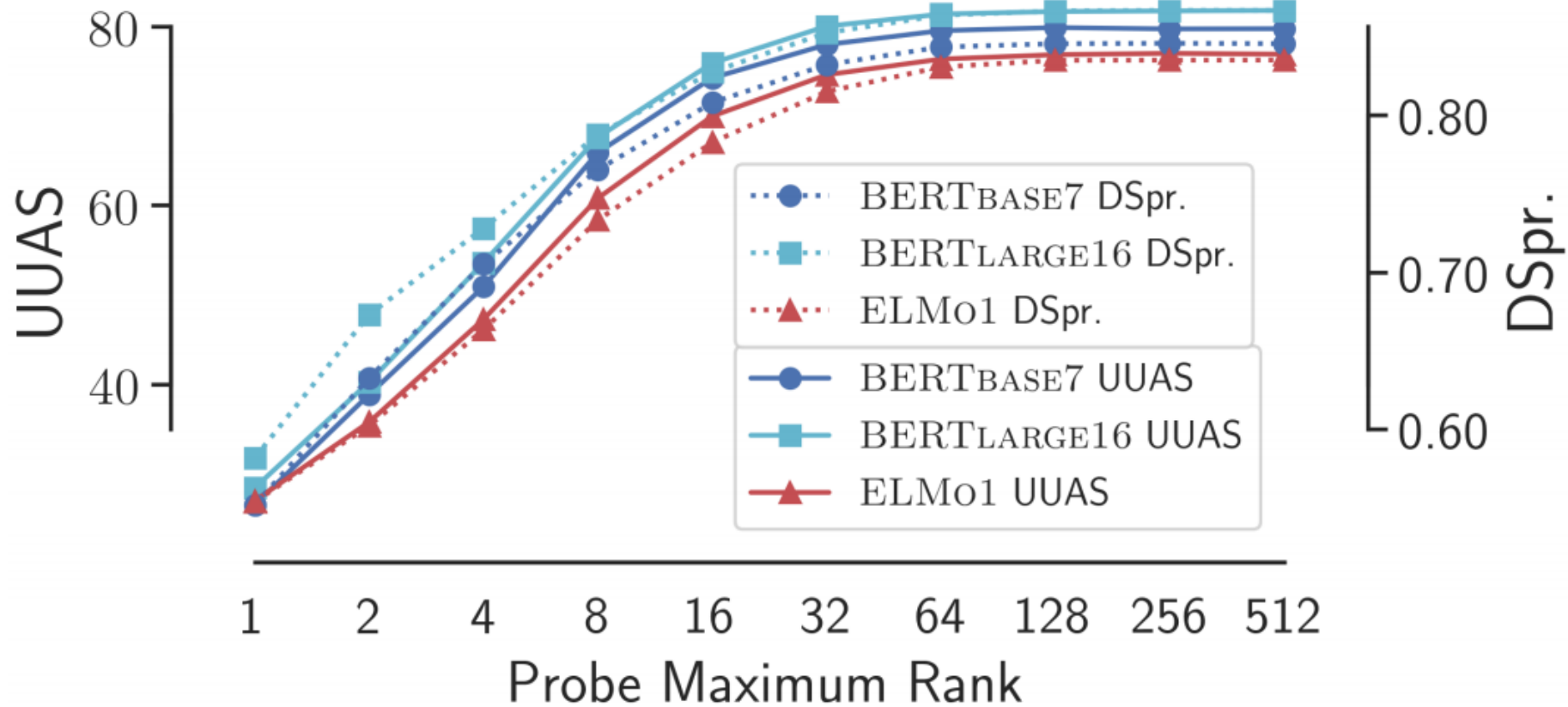


<Parse Distance of Predicted & Gold Tree (from BERT_{LARGE16})>

Experiments

-Result

<Result>



<Parse Distance Tree Reconstruction Accuracy on Various Probe Maximum Rank(k)>

Future Works

Future Work

-More Things to do in the Future

<More Things to do in the Future>

- **Why Does Squared L2 Distance Reconstruct Dependency Parse Tree?**
- **Representation of Language Model Does Not Use Full Dimension to Contain Dependency Relation (Syntactic Information) of Text.**
- **Then, What More Information is Contained in Language Model Embedding?**

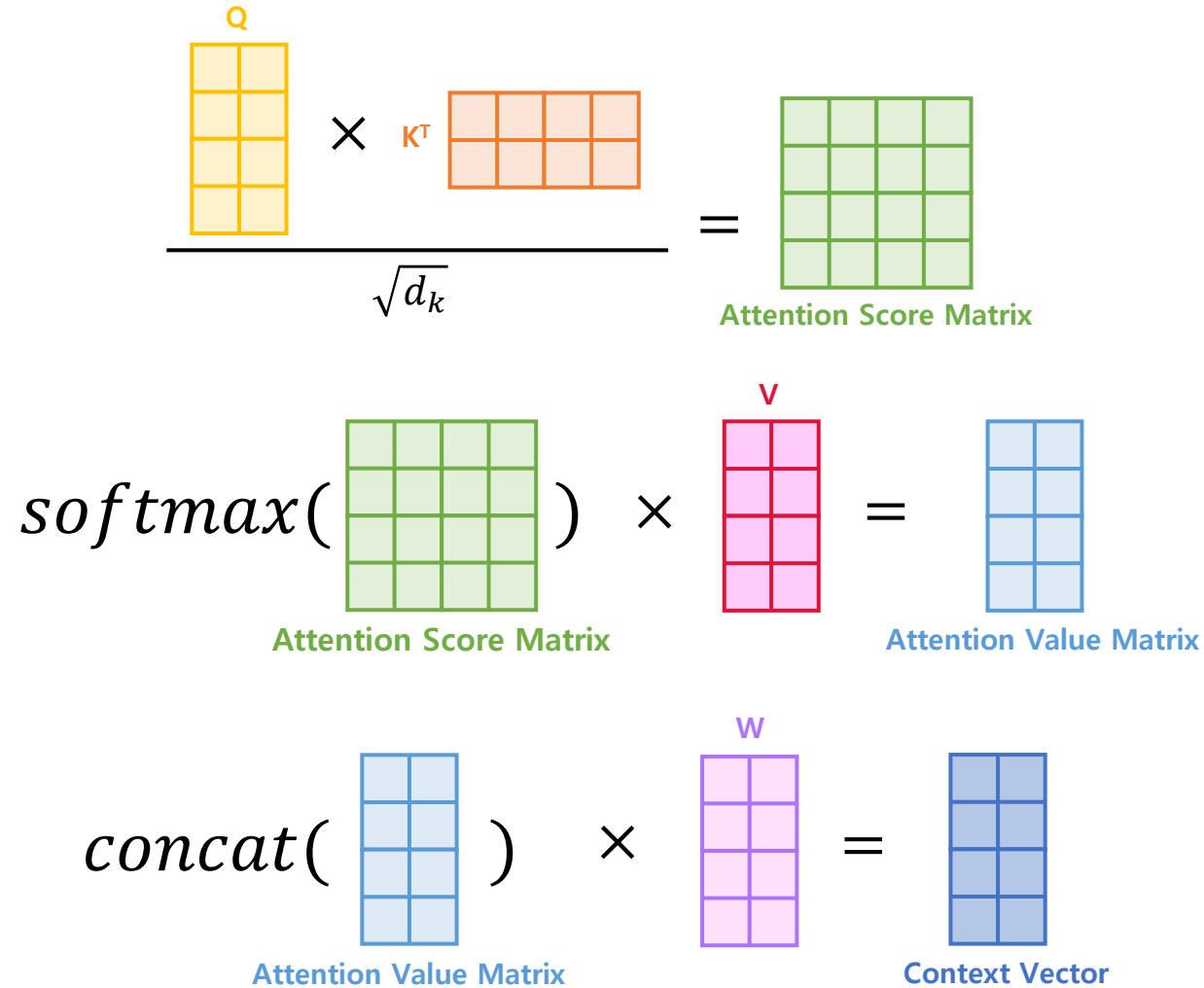
Visualizing and Measuring the Geometry of BERT

Coenen et al., 2019, NIPS

Geometry of Syntax

- **Attention Probe**
- **Geometry of Parse Tree Embedding**

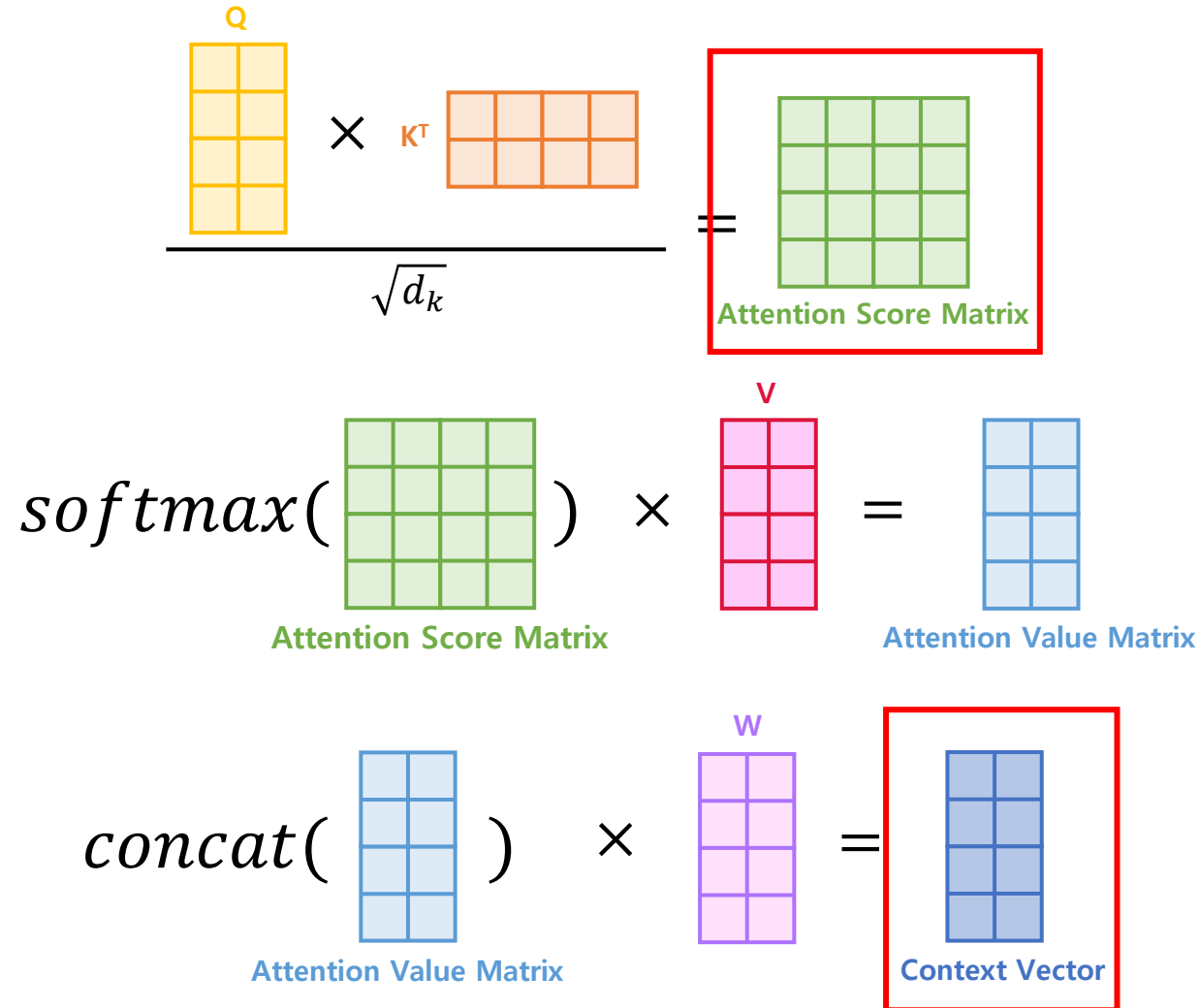
<Multi-head Self Attention>



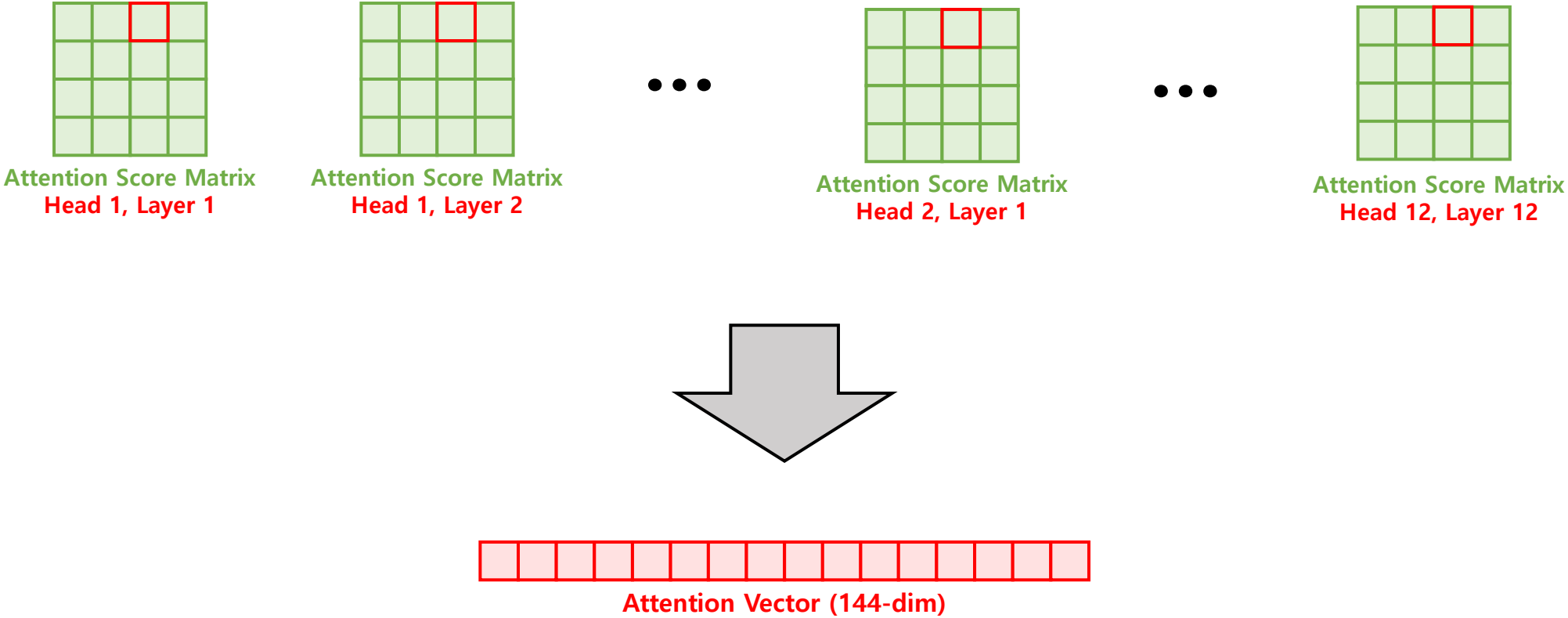
Geometry of Syntax

-Attention Probe

<Attention Probe>

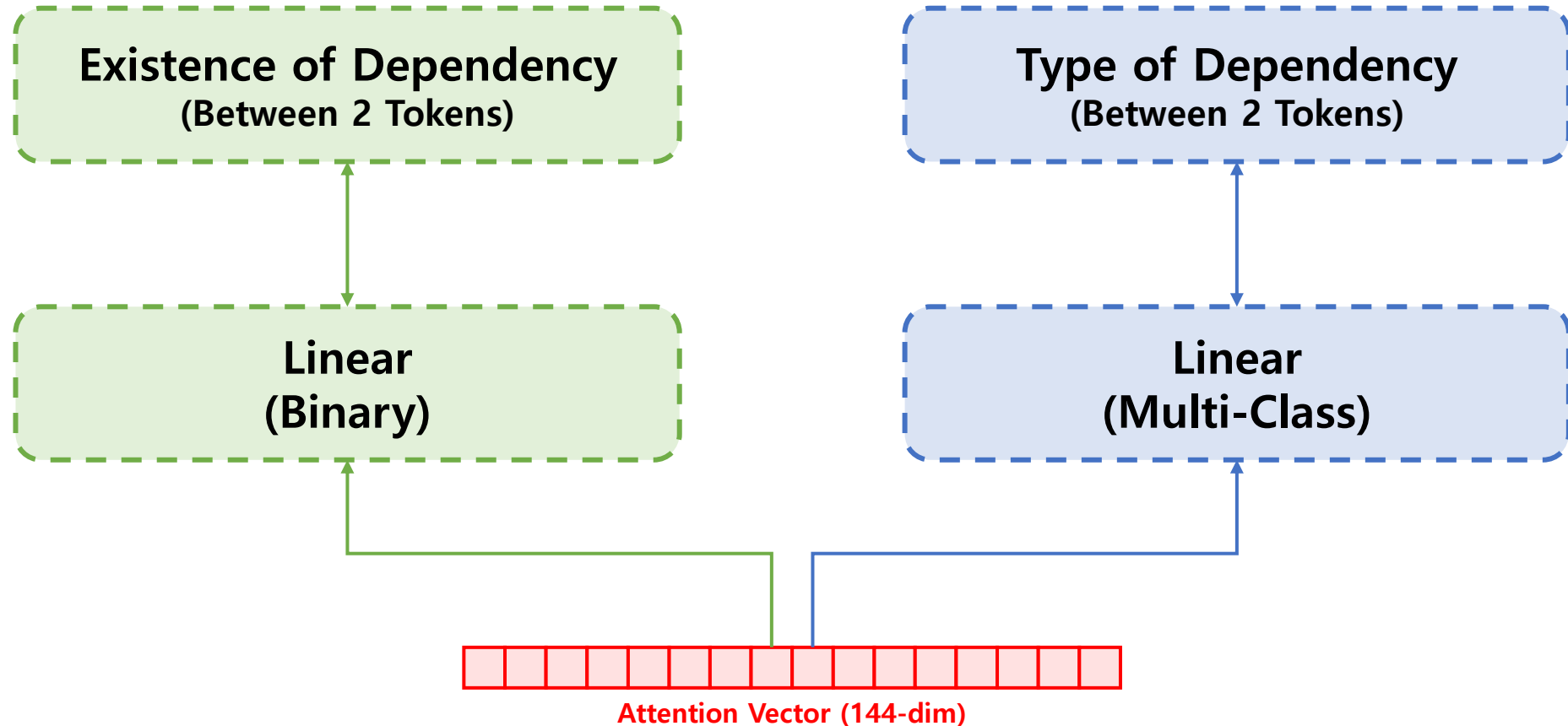


<Attention Probe>



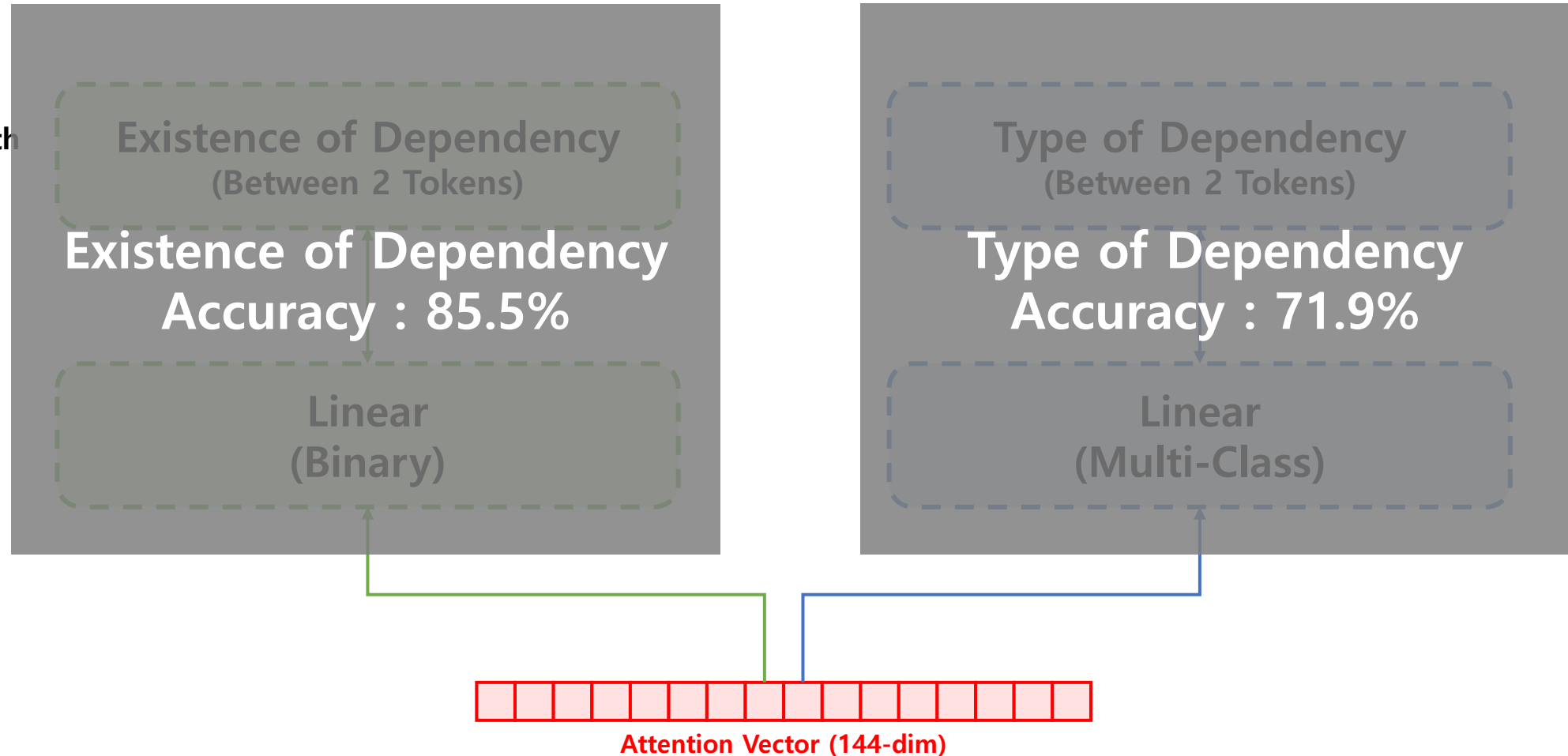
<Attention Probe>

Ground Truth
of PBE



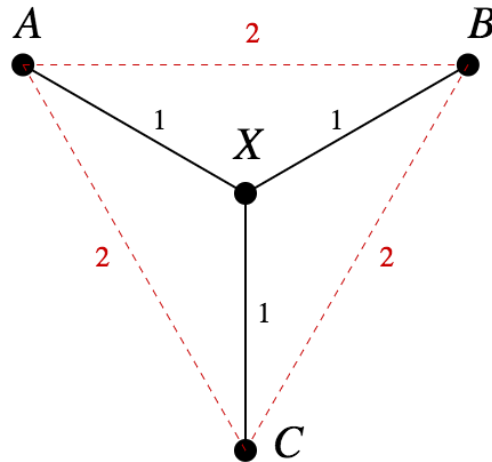
<Attention Probe>

Ground Truth
of PBE



<Pythagorean Embedding>

“One cannot generally embed a tree, with its tree metric d , isometrically into Euclidean space”



$$d(A, B) = d(A, X) + d(X, B)$$

$\therefore A, X \text{ and } B \text{ is collinear}$

$$d(A, C) = d(A, X) + d(X, C)$$

$\therefore A, X \text{ and } C \text{ is collinear}$

$\therefore B = C : \text{contradiction}$

Geometry of Syntax

-Geometry of Parse Tree Embedding

<Pythagorean Embedding>

Tree $T: t_0, \dots, t_{n-1}$, where t_0 : root node

$\{e_1, \dots, e_{n-1}\}$: orthogonal unit basis vectors for \mathbb{R}^{n-1}

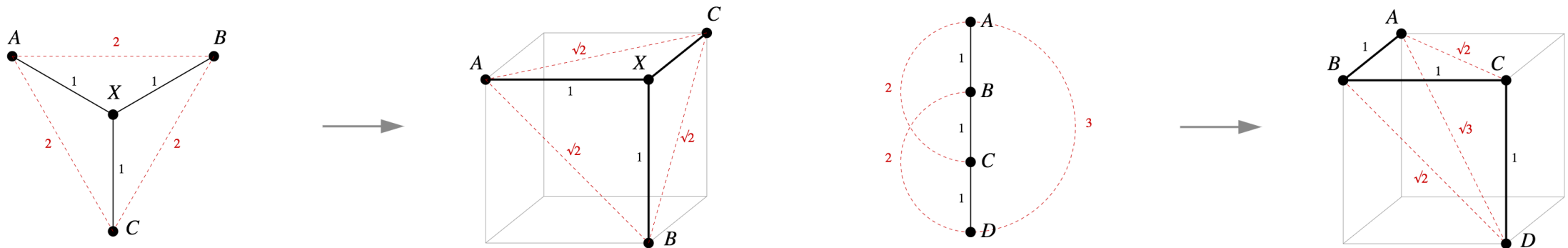
$$f: T \rightarrow \mathbb{R}^{n-1}$$

$$f(t_0) = \mathbf{0}$$

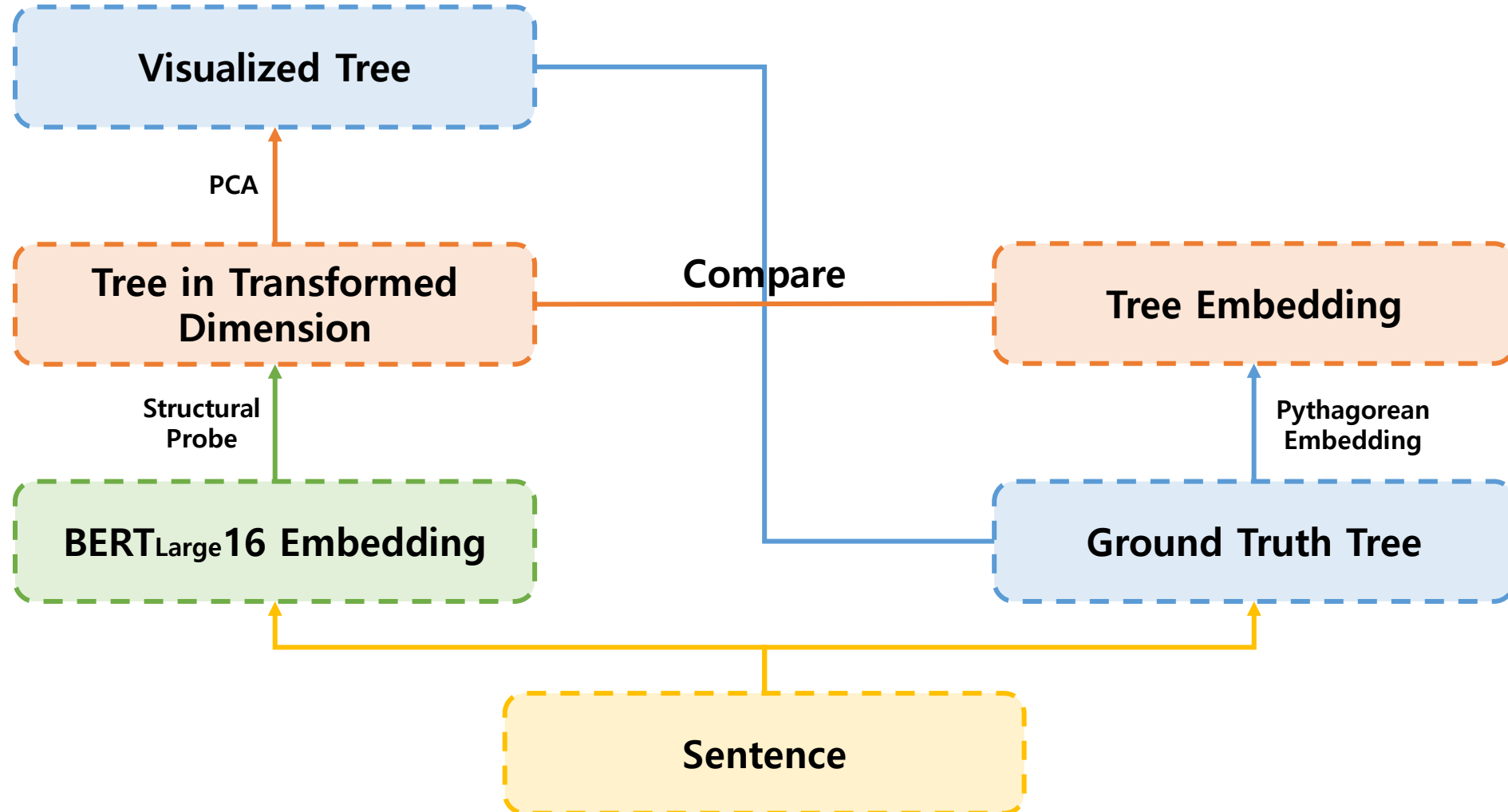
$$f(t_i) = e_i + f(\text{parent}(t_i))$$

$$\|f(x) - f(y)\|^2 = m = d(x, y)$$

where, n : number of nodes, $m = d(x, y)$: tree distance



<Visualization of Parse Tree Embedding>

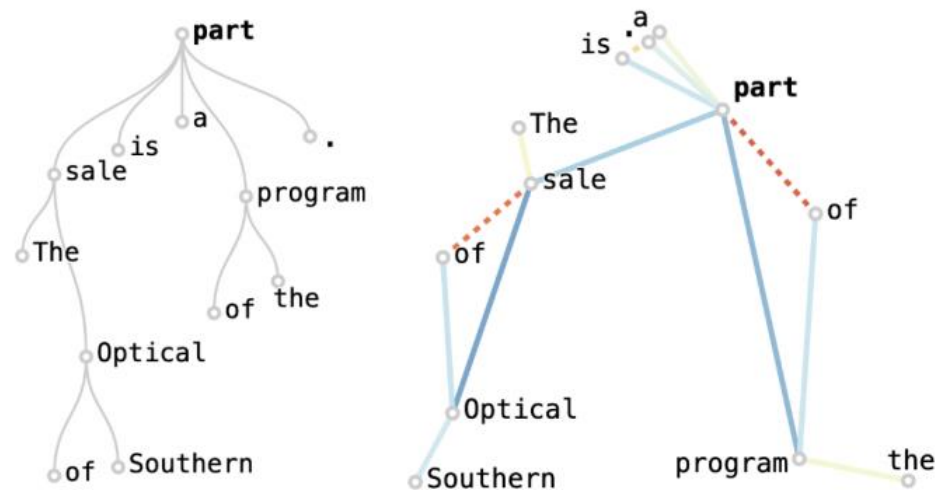


Geometry of Syntax

-Visualization of Parse Tree Embedding

<Visualization of Parse Tree Embedding>

"The sale of Southern Optical is a part of the program."

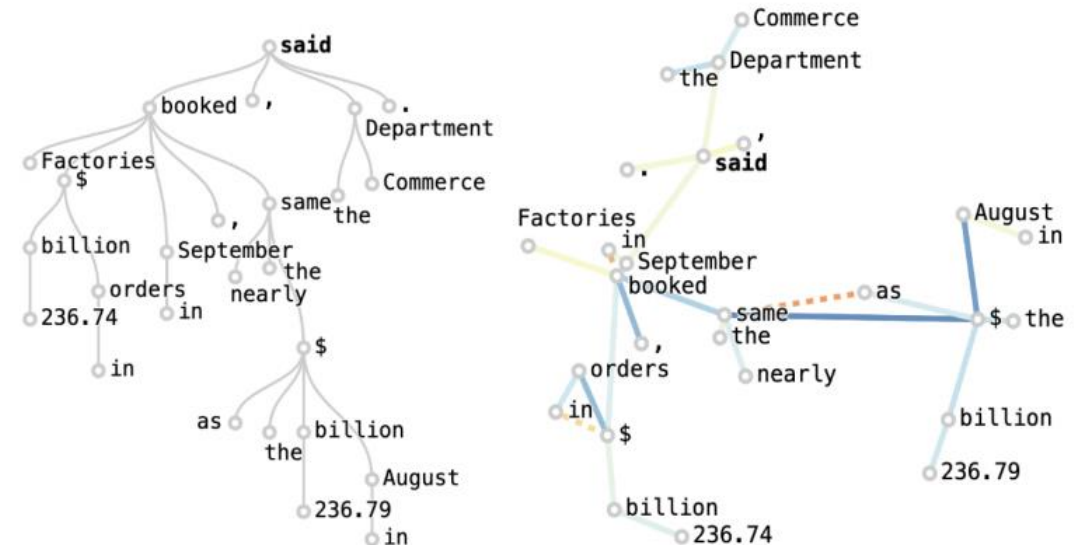


Ratio between d^2 and tree distance



— Ground truth dependency
.... No ground truth dependency, $d^2 < 1.5$

"Factories booked \$236.74 billion in orders in September, nearly the same as the \$236.79 billion in August, the Commerce Department said."



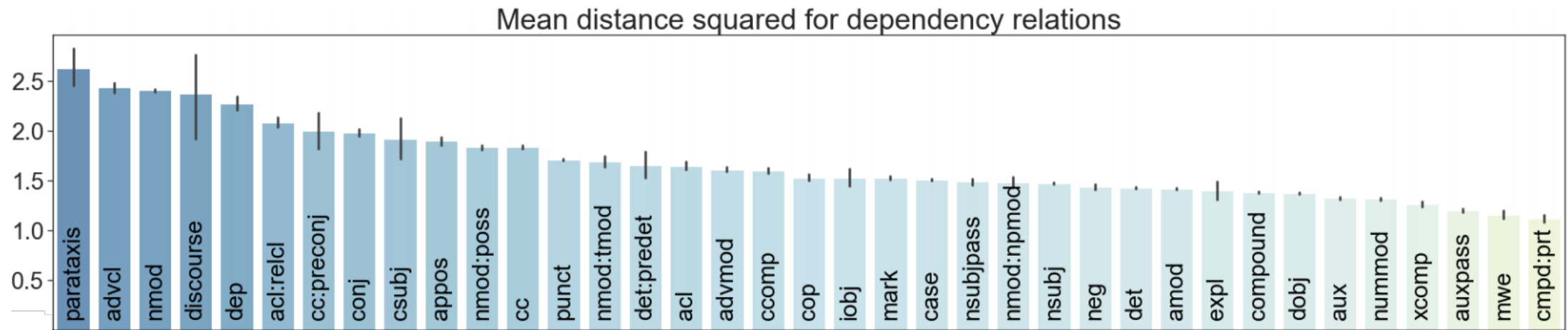
<Visualizations of Tree Embedding>

(Left – Parse Tree, Right – PCA Projection of Context Embedding)

Geometry of Syntax

-Visualization of Parse Tree Embedding

<Visualization of Parse Tree Embedding>



<The Average Squared Edge Length Between Two Words with a Given Dependency>

“BERT’s syntactic representation has an additional quantitative aspect beyond traditional dependency grammar”

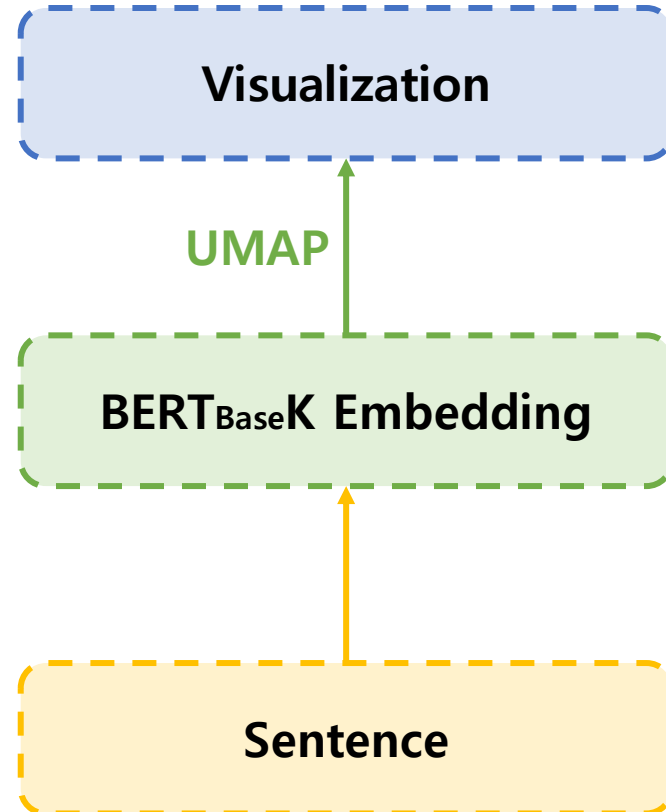
Geometry of Word Senses

- Visualization of Word Senses
- Measurement of Word Sens Disambiguation
- Embedding Distance and Context

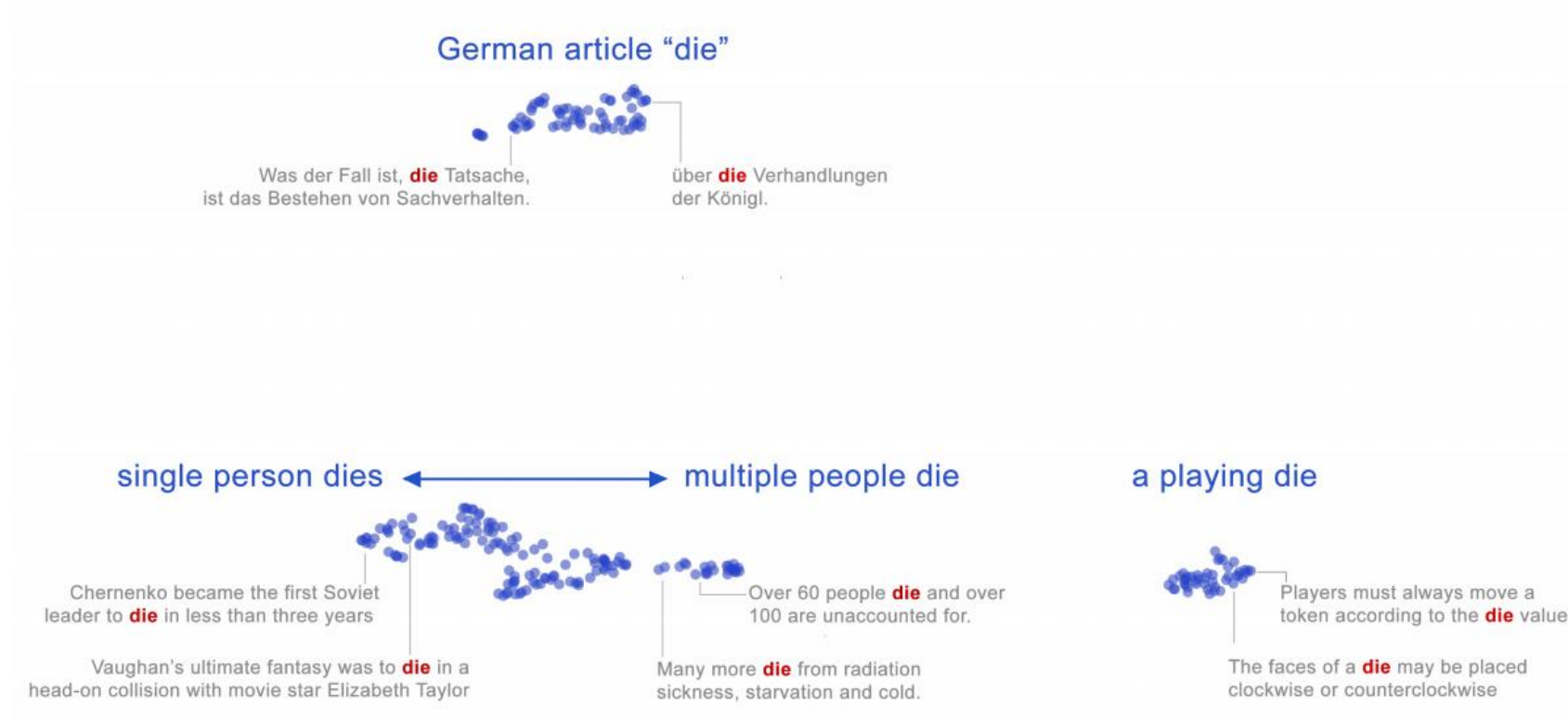
Geometry of Word Senses

-Visualization of Word Senses

<Visualization of Word Senses>



<Visualization of Word Senses>



<Embeddings for the Word "die" in Different Contexts>

Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Word Sense Disambiguation>

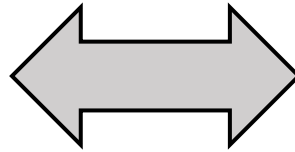
"Give me an **account** of what you saw."

"The problem is important on this **account**."

"You must give in my **account** once a month."

"In his **account** it was very excellent."

<Sentence>



Sense 1: Explanation

Sense 2: Reason

Sense 3: Bank account

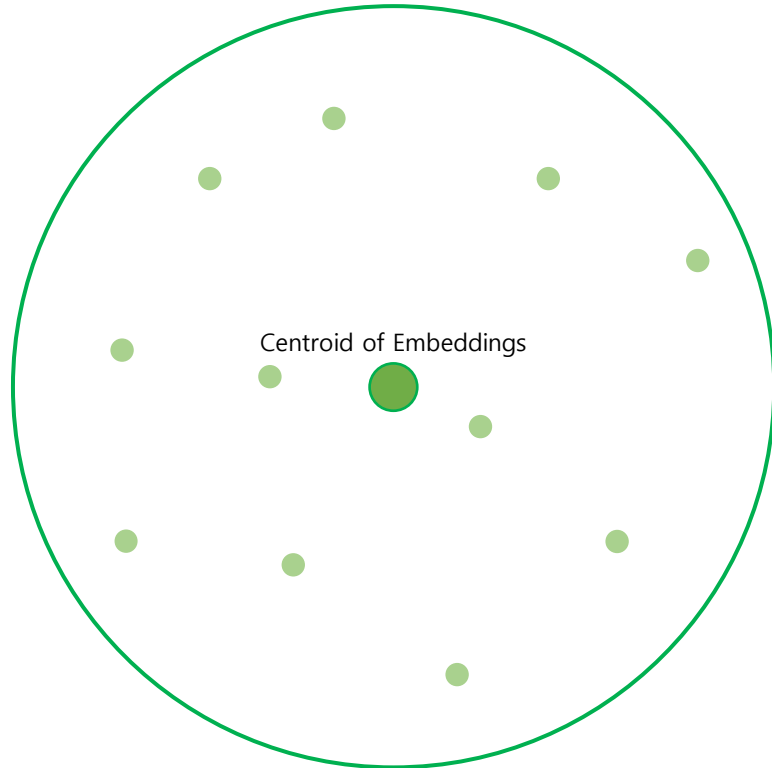
Sense 4: Evaluation

<Senses of all Tokens in Given Sentence>

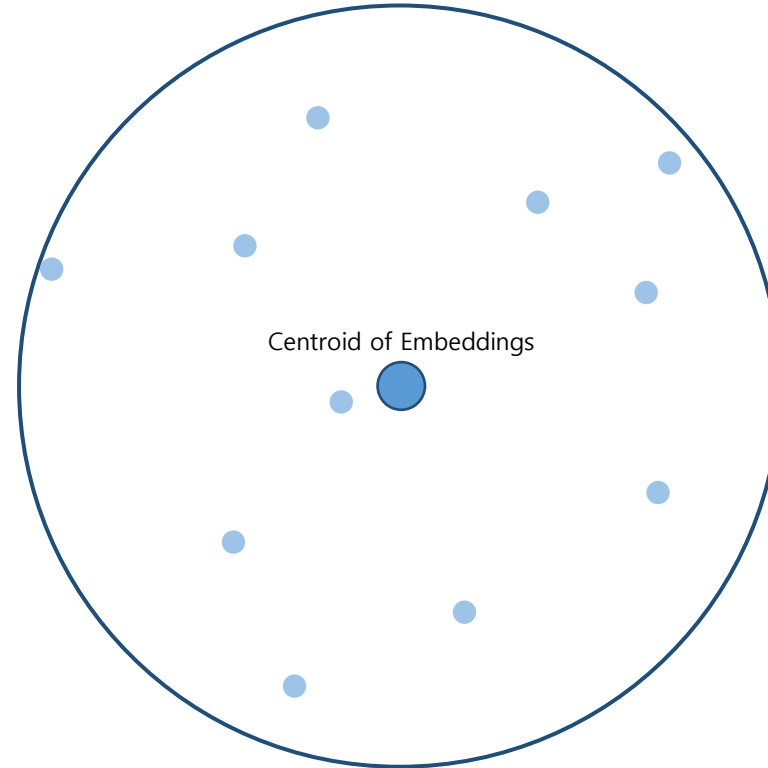
Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Nearest-neighbor Classification>



<Embeddings of a Word with Sense 1>

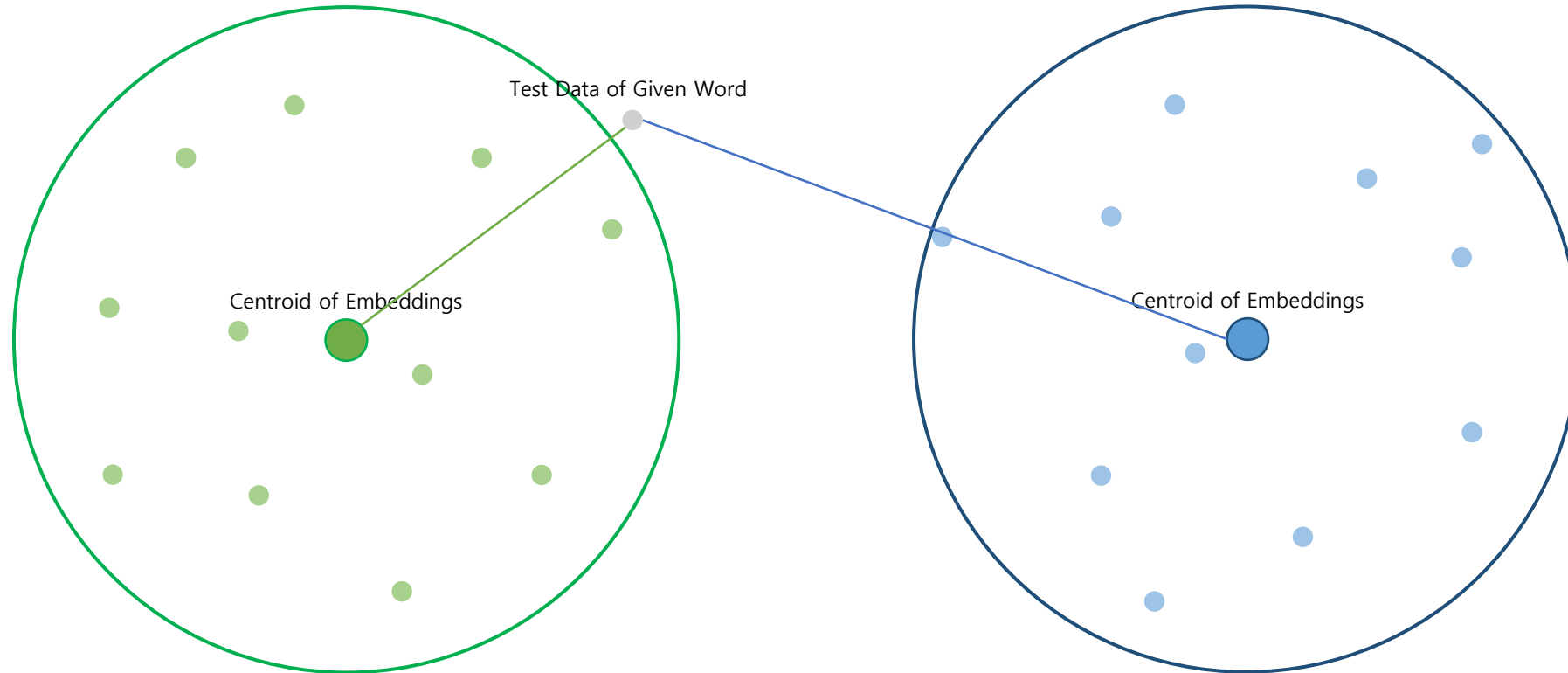


<Embeddings of a Word with Sense 2>

Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Nearest-neighbor Classification>



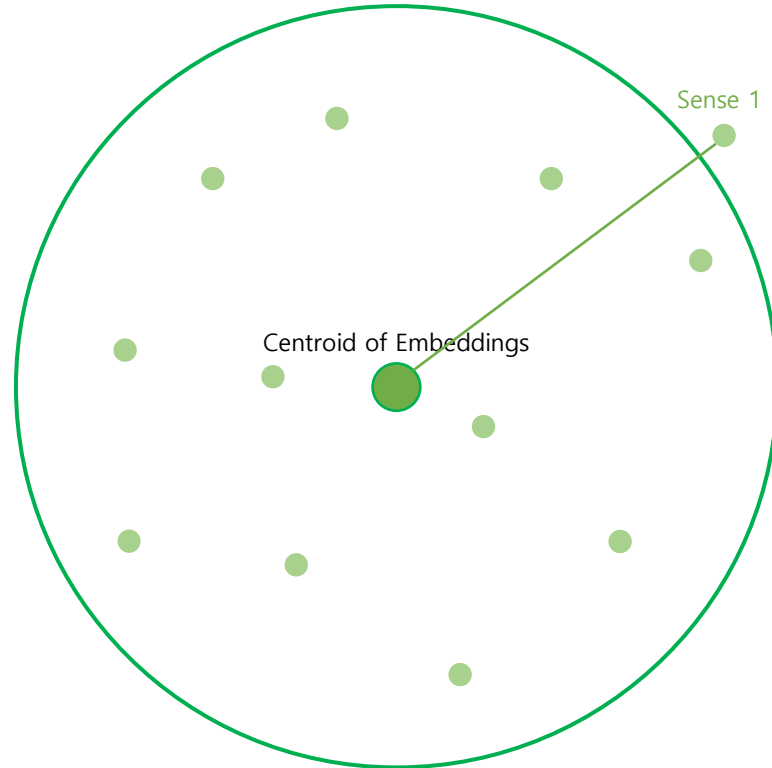
<Embeddings of a Word with Sense 1>

<Embeddings of a Word with Sense 2>

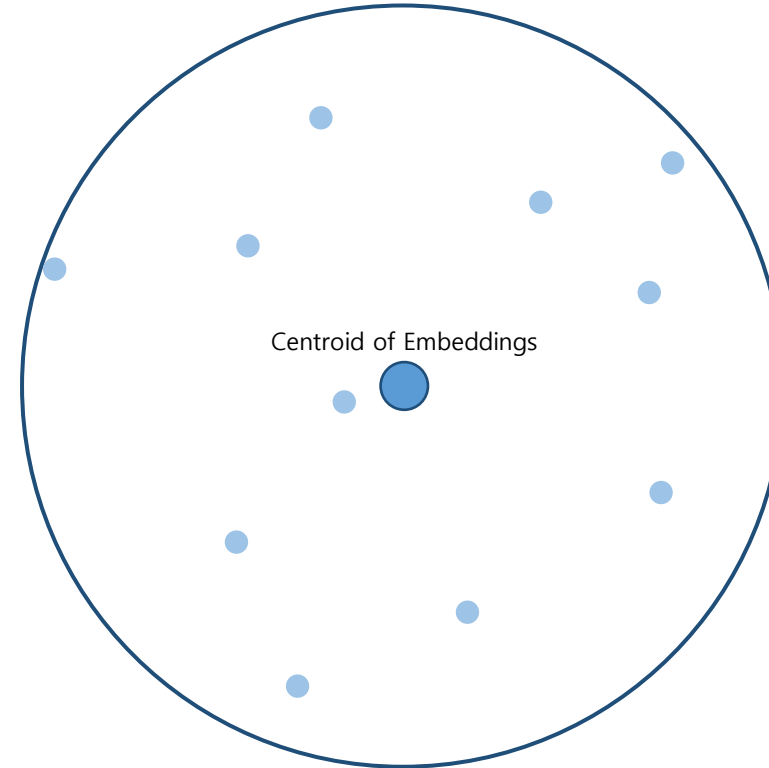
Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Nearest-neighbor Classification>



<Embeddings of a Word with Sense 1>



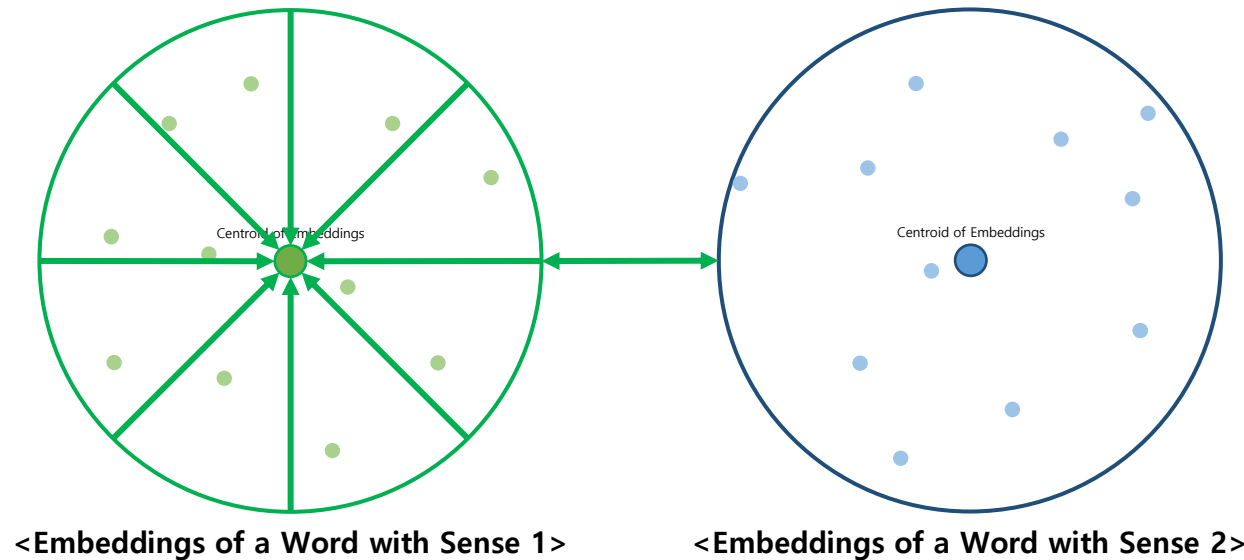
<Embeddings of a Word with Sense 2>

Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Structural Probe for Word Sense Disambiguation>

$\min_B(\text{average cosine similarity of same sense} - \text{average cosine similarity of different sense})$



<Result>

Corpus	Method	F1 Score
SemCor	IMS	68.4
	IMS+emb	69.1
	IMS _{-s} +emb	69.6
	Context2Vec	69.0
	MFS	64.8

<Raganato et al. 2017>

Corpus	Method	F1 Score
SemCor	Baseline (most frequent sense)	64.8
	ELMo	70.1
	BERT	71.1
	BERT (w/probe)	71.5

<This Paper>

<Result>

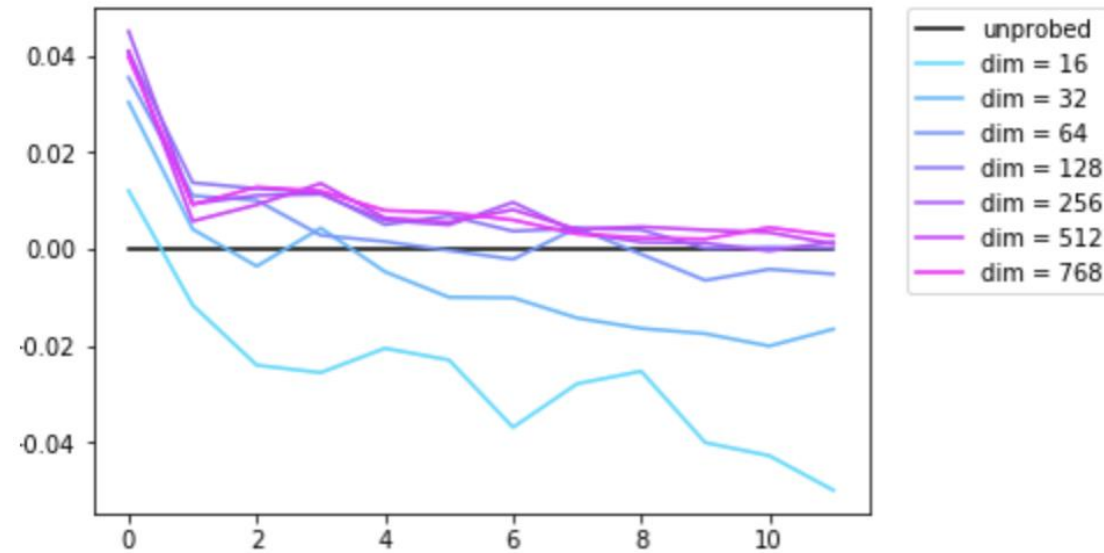
m	Trained probe	Random probe
768 (full)	71.26	70.74
512	71.52	70.51
256	71.29	69.92
128	71.21	69.56
64	70.19	68.00
32	68.01	64.62
16	65.34	61.01

<Semantic Probe % Accuracy on Final-layer BERT-base>

Geometry of Word Senses

-Measurement of Word Sense Disambiguation

<Result>



<Change in Classification Accuracy by Layer of Different Probe Dimensionalities>

<Embedding Distance and Context>

Sentence A: "He thereupon went to London and spent the winter talking to men of wealth."

went: to move from one place to another (**Sense 1**)

Sentence B: "He went prone on his stomach, the better to pursue his examination."

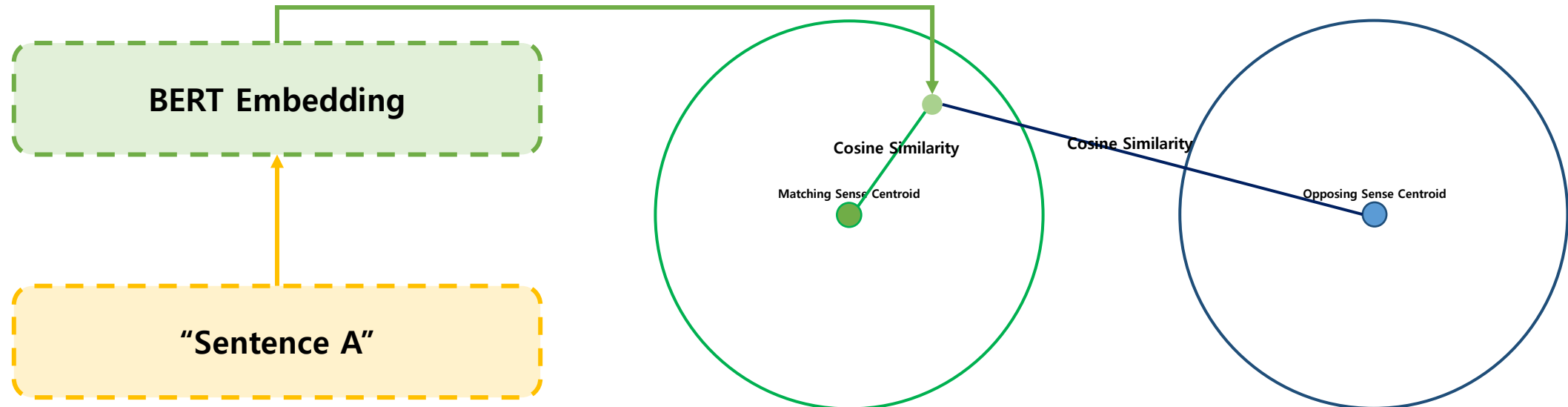
went: to enter into a specified state. (**Sense 2**)

Sense 1: matching sense of Sentence A

Sense 2: opposing sense of Sentence A

<Embedding Distance and Context>

$$\text{individual similarity ratio} = \frac{\text{matching sense similarity}}{\text{opposing sense similarity}}$$

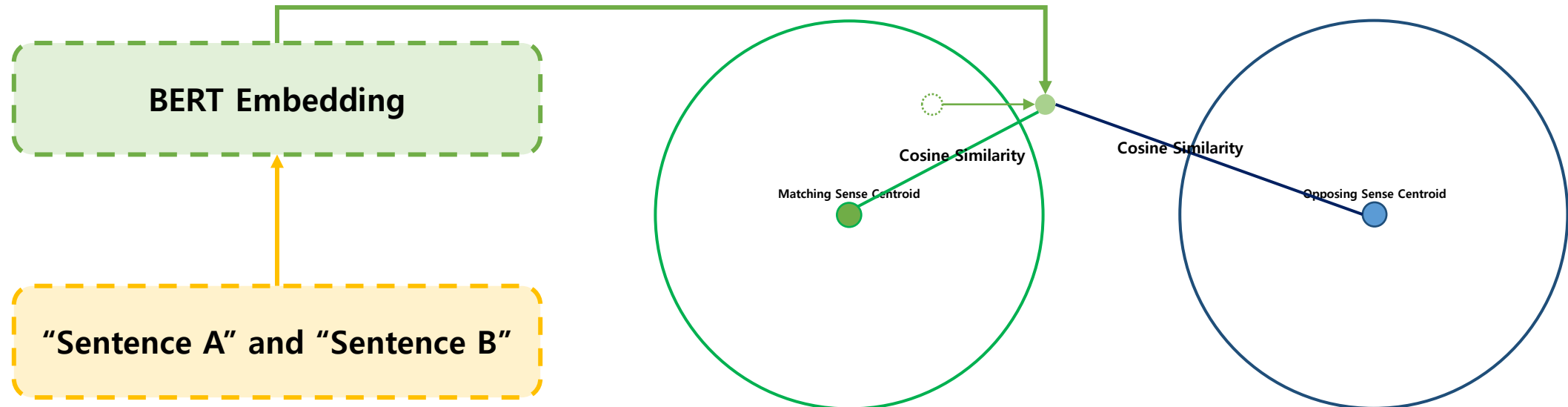


Geometry of Word Senses

-Embedding Distance and Context

<Embedding Distance and Context>

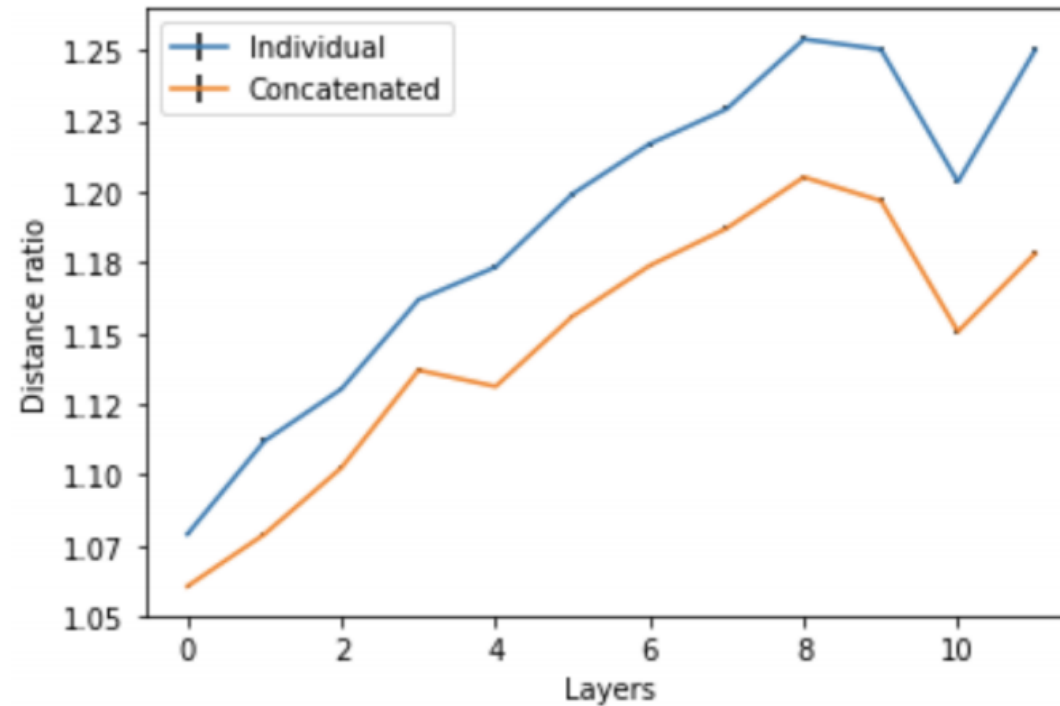
$$\text{concatenated similarity ratio} = \frac{\text{matching sense similarity}}{\text{opposing sense similarity}}$$



Geometry of Word Senses

-Embedding Distance and Context

<Result>



<Average Similarity Ratio: Senses A vs. B>

Conclusion

<Conclusion>

- **There are Subspaces in Language Model Representation that Contain Syntactic and Semantic Information Respectively**
- **There are Limitations of Attention-based Model: Tokens do not Respect Semantic Boundaries, But Absorb Meaning from all Neighbors**

Any Questions?

Thank You