

Paper Review

Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision

Tan and Bansal, 2020, EMNLP

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

Introduction

- Human Language Acquisition

Introduction

-What This Seminar Does Not Cover

<What This Seminar Does Not Cover>

- **Details of BERT**

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

Introduction

-Recap: Previous Seminar

From [Paper Review] Climbing towards NLU (Seungwan Seo, 2020)



<Recap: Previous Seminar>

Human language acquisition

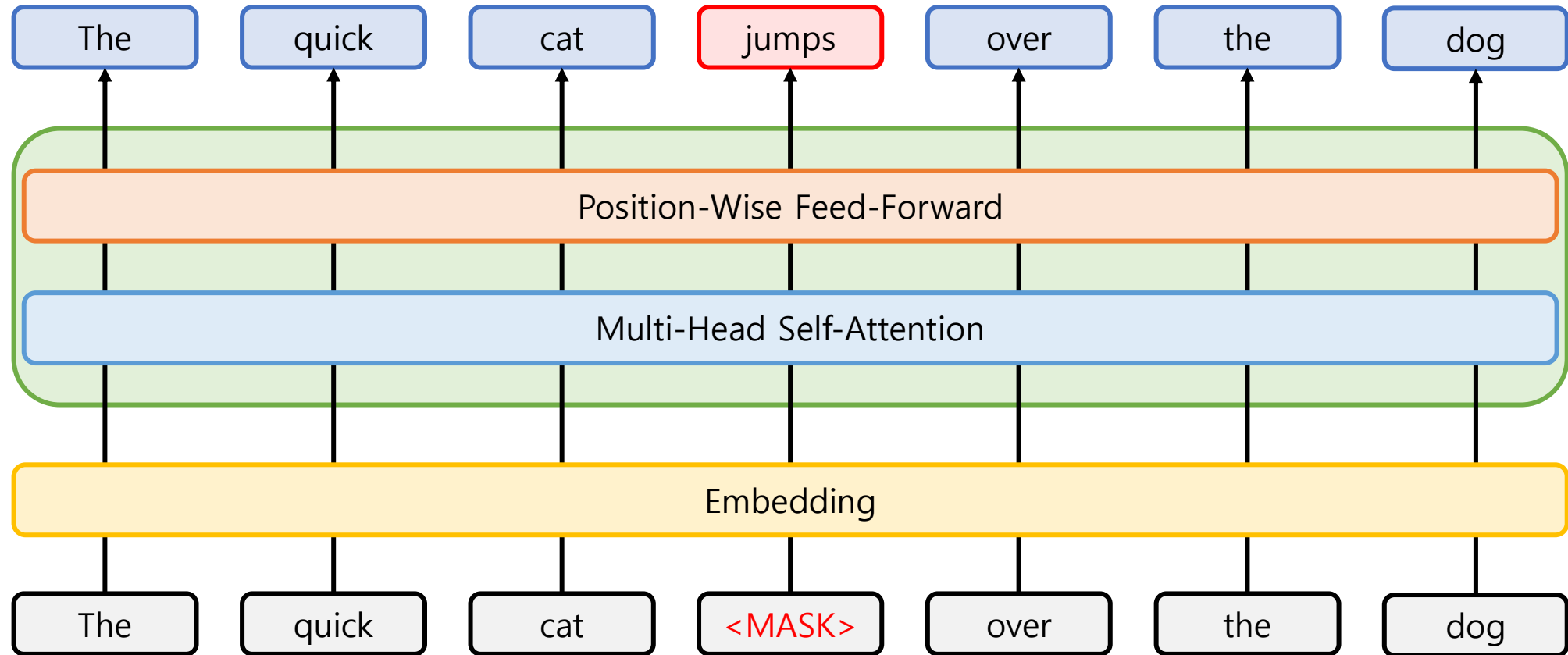


[사람은 결국 언어의 형태(표현) 뿐만 아니라,
다양한 요소의 상호작용을 통해 언어를 습득함]

Introduction

-Human Language Acquisition

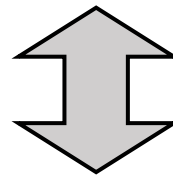
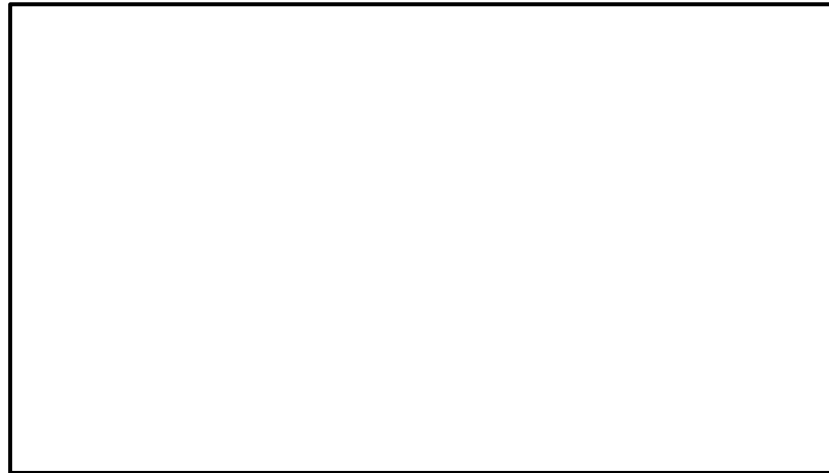
<Model Language Acquisition>



Introduction

-Human Language Acquisition

<Human Language Acquisition>

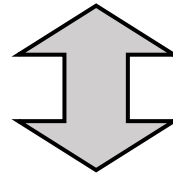


"The quick cat jumps over the dog."

Introduction

-Human Language Acquisition

<Human Language Acquisition>



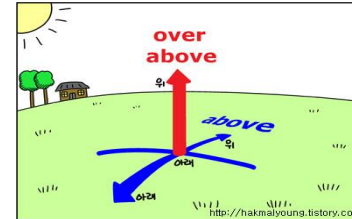
"The quick cat jumps over the dog."

Introduction

-Human Language Acquisition

<Human Language Acquisition>

the



the

The

quick

cat

jumps

over

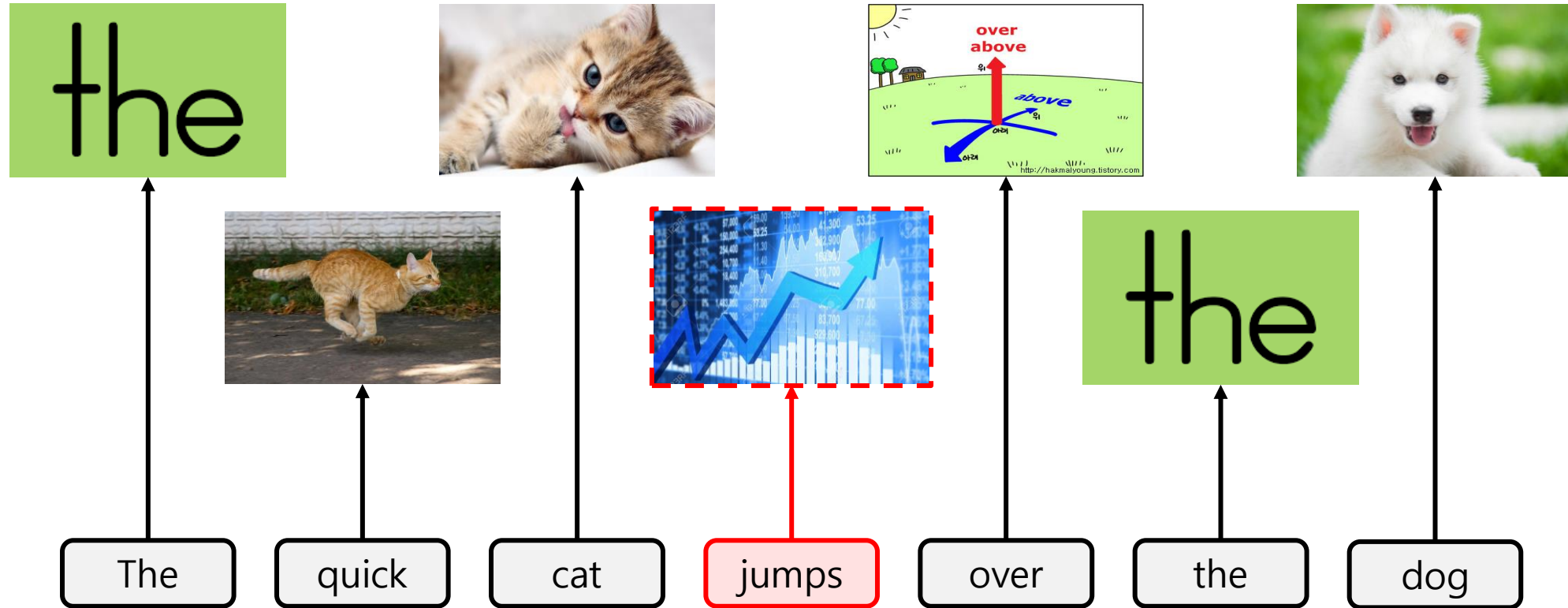
the

dog

Introduction

-Human Language Acquisition

<Human Language Acquisition>

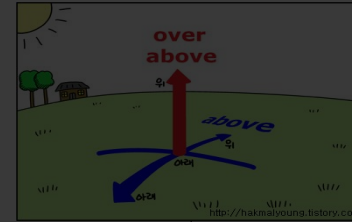
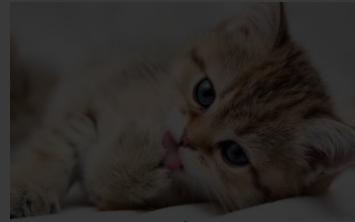


Introduction

-Human Language Acquisition

<Human Language Acquisition>

Can the Language Model Learn Contextual Information from the External Visual World?



The

quick

cat

jumps

over

the

dog

Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision

Tan and Bansal, 2020, EMNLP

Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision

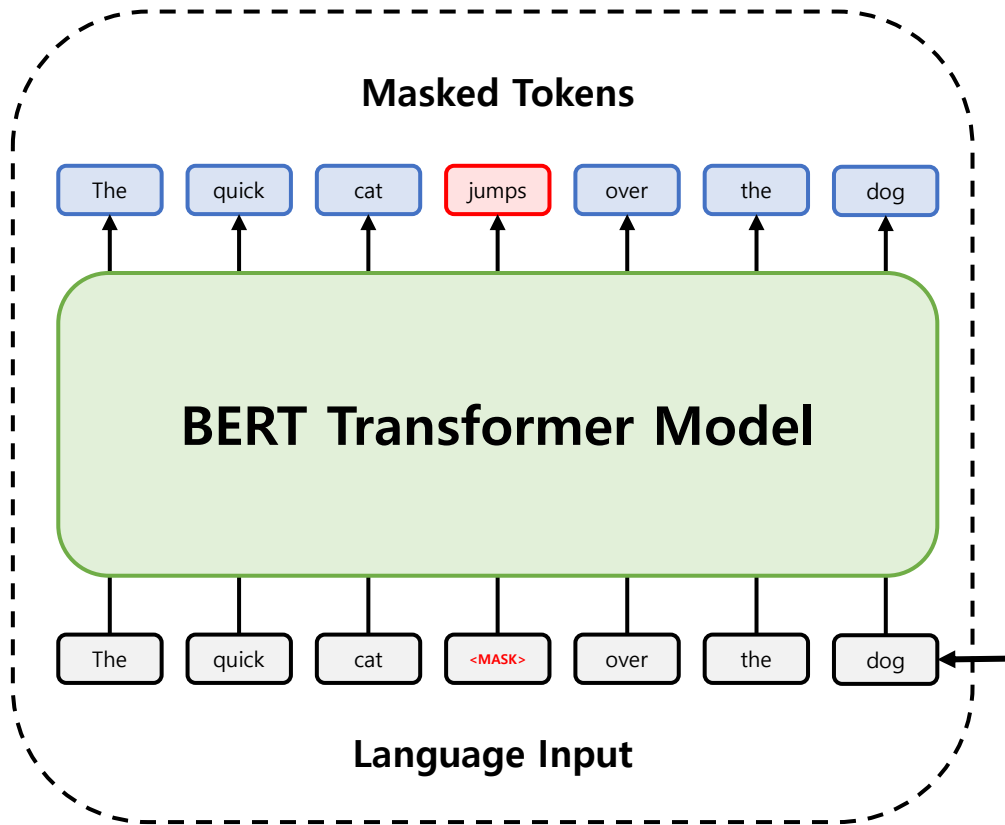
Tan and Bansal, 2020, EMNLP

Model

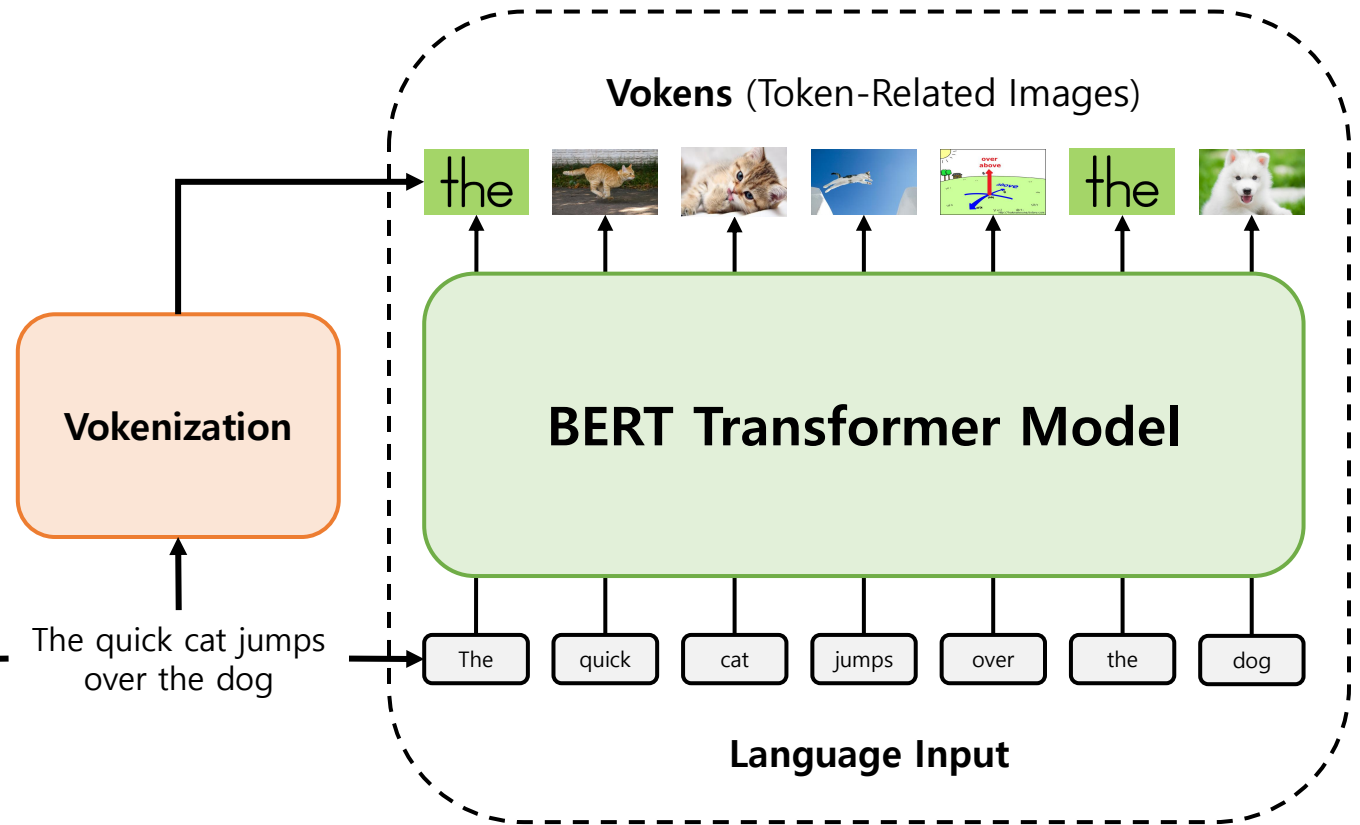
- The Voken-Classification Task
- Visually-Supervised Language Model
- Two Challenges in Creating Vokens

<Visually-Supervised Language Model>

Masked Language Model



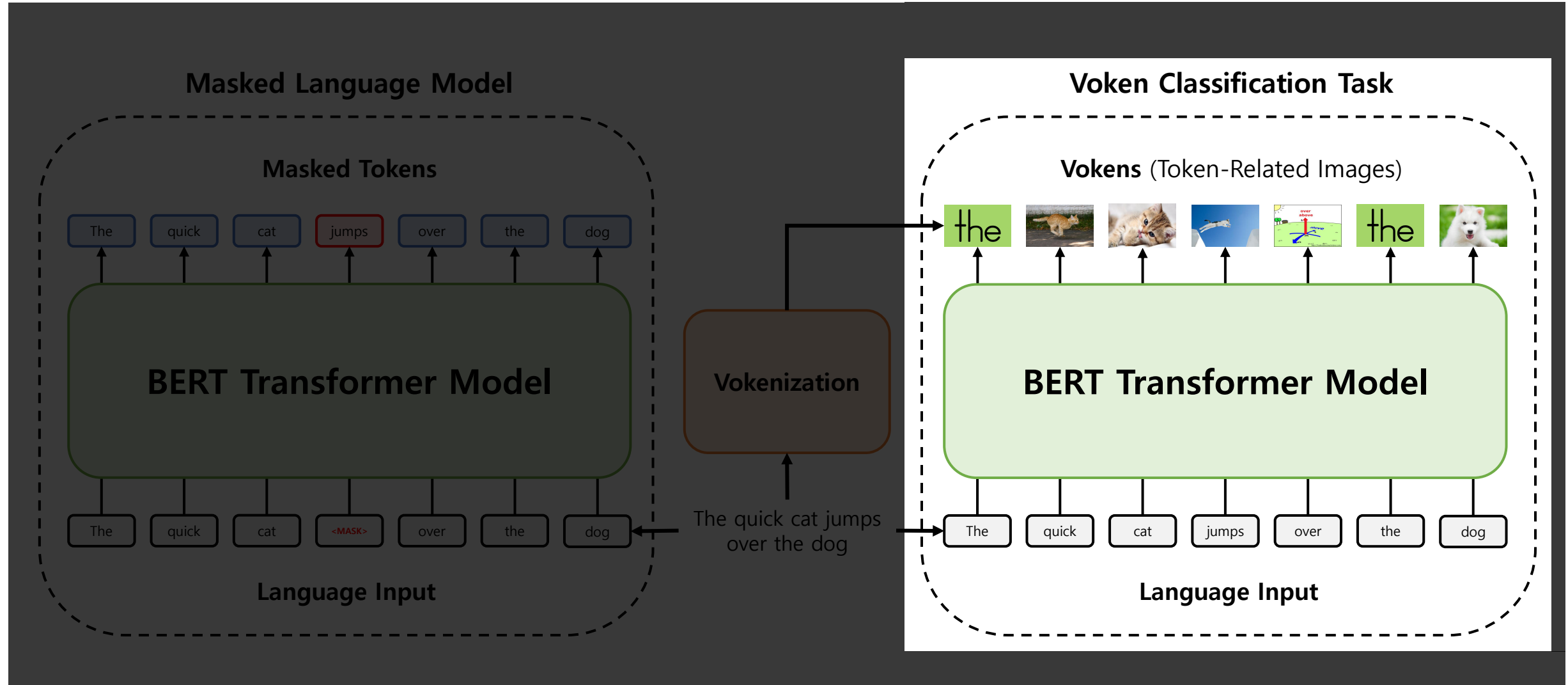
Voken Classification Task



Model

-The Voken Classification Task

<The Voken Classification Task>



Model

-The Voken Classification Task

<The Voken Classification Task>

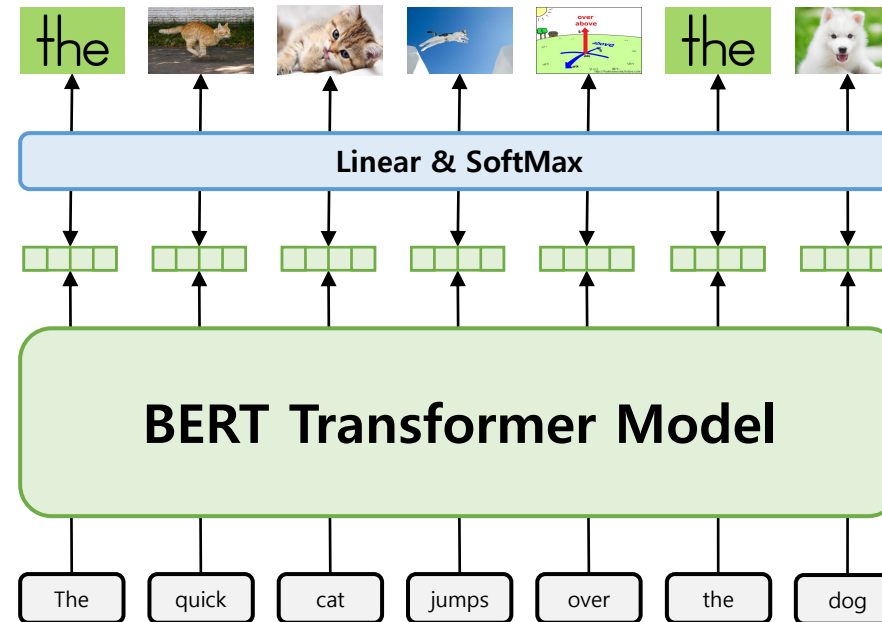
$\{h_i\}$: Language Model Representations

$s = \{w_i\}$, where S : Sentence, w_i : Tokens

$h_1, h_2, \dots, h_l = LM(w_1, w_2, \dots, w_l)$

$p_i(v|s) = \text{softmax}_v\{Wh_i + b\}$

$$\mathcal{L}_{VOKEN-CLS}(s) = - \sum_{i=1}^l \log p_i(v(w_i ; s) | s)$$



Model

-The Voken Classification Task

<The Voken Classification Task>

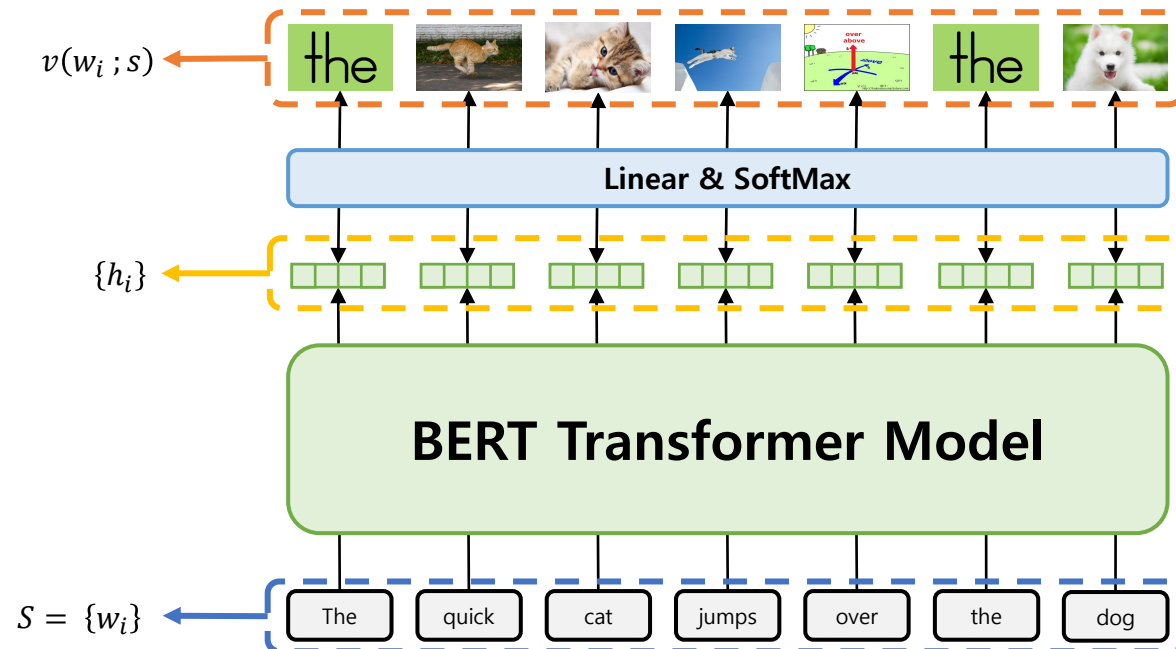
$\{h_i\}$: Language Model Representations

$s = \{w_i\}$, where S : Sentence, w_i : Tokens

$h_1, h_2, \dots, h_l = LM(w_1, w_2, \dots, w_l)$

$p_i(v|s) = \text{softmax}_v\{Wh_i + b\}$

$$\mathcal{L}_{VOKEN-CLS}(s) = - \sum_{i=1}^l \log p_i(v(w_i; s) | s)$$



Model

-Visually-Supervised Language Model

<Visually-Supervised Language Model>

Masked Language Model

Masked Tokens

The quick cat jumps over the dog

BERT Transformer Model

The quick cat <MASK> over the dog

Language Input

Vokenization

The quick cat jumps
over the dog

Voken Classification Task

Vokens (Token-Related Images)

the the the the the the the

BERT Transformer Model

The quick cat jumps over the dog

Language Input

Model

-Visually-Supervised Language Model

<Visually-Supervised Language Model>

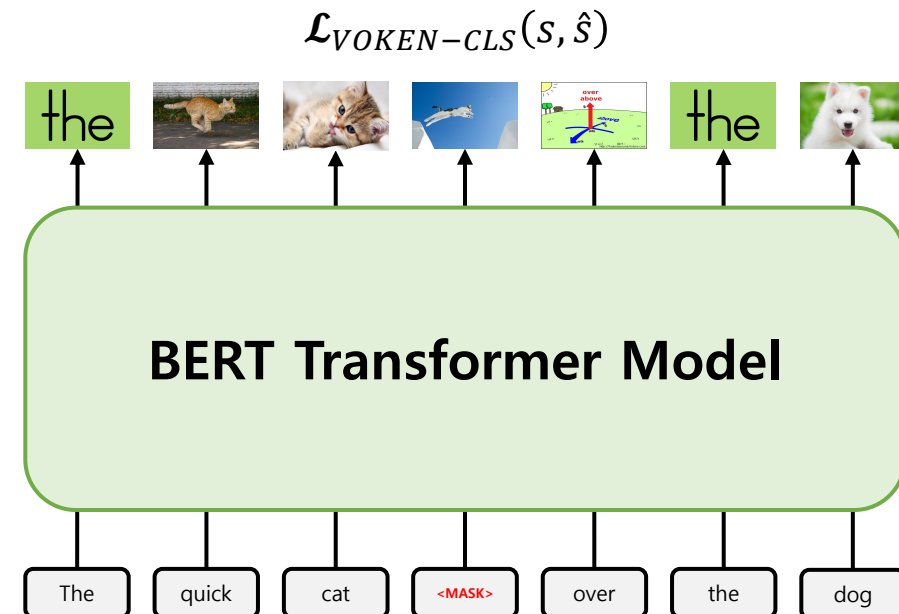
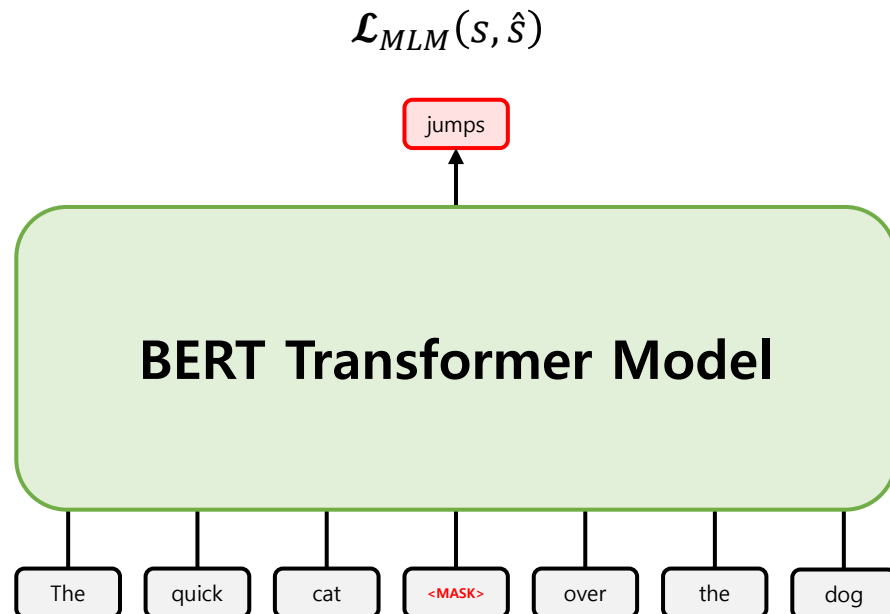
s, \hat{s} : Set of Tokens, Masked Tokens

$s \setminus \hat{s}$: Unmasked Tokens

$$\mathcal{L}_{MLM}(s, \hat{s}) = - \sum_{w_i \in \hat{s}} \log q_i(w_i | s \setminus \hat{s})$$

$$\mathcal{L}_{VOKEN-CLS}(s, \hat{s}) = - \sum_{w_i \in s} \log p_i(v(w_i; s) | s \setminus \hat{s})$$

$$\mathcal{L}_{VLM}(s, \hat{s}) = \mathcal{L}_{VOKEN-CLS}(s, \hat{s}) + \lambda \mathcal{L}_{MLM}(s, \hat{s}) \quad (1)$$



Model

-Two Challenges in Creating Vokens

<Two Challenges in Creating Vokens>

It is Difficult to Apply **Image Captioning** to Large Corpus for the Following Two Reasons:

(1) Different Distributions between Grounded Language and Other Natural Language Corpora



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

Average Sentence Length = 11.8 words

<MS COCO 2015 Image Captioning Task>

BERT (language model)

From Wikipedia, the free encyclopedia

Bidirectional Encoder Representations from Transformers (BERT) is a [Transformer](#)-based [machine learning](#) technique for [natural language processing](#) (NLP) pre-training developed by [Google](#). BERT was created and published in 2018 by Jacob Devlin and his colleagues from Google.^{[1][2]} As of 2019, Google has been leveraging BERT to better understand user searches.^[3]

The original English-language BERT model comes with two pre-trained general types:^[1] (1) the BERT_{BASE} model, a 12-layer, 768-hidden, 12-heads, 110M parameter neural network architecture, and (2) the BERT_{LARGE} model, a 24-layer, 1024-hidden, 16-heads, 340M parameter neural network architecture; both of which were trained on the [BooksCorpus](#)^[4] with 800M words, and a version of the [English Wikipedia](#) with 2,500M words.

Average Sentence Length = 24.1 words

<Wikipedia Database>

Model

-Two Challenges in Creating Vokens

<Two Challenges in Creating Vokens>

It is Difficult to Apply **Image Captioning** to Large Corpus for the Following Two Reasons:

(2) Low Grounding Ratio in Natural Language



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

Num of Sentences : 410K

Num of Images : 82K

<MS COCO 2015 Image Captioning Task>

BERT (language model)

From Wikipedia, the free encyclopedia

Bidirectional Encoder Representations from Transformers (BERT) is a [Transformer](#)-based [machine learning](#) technique for [natural language processing](#) (NLP) pre-training developed by [Google](#). BERT was created and published in 2018 by Jacob Devlin and his colleagues from Google.^{[1][2]} As of 2019, Google has been leveraging BERT to better understand user searches.^[3]

The original English-language BERT model comes with two pre-trained general types:^[1] (1) the BERT_{BASE} model, a 12-layer, 768-hidden, 12-heads, 110M parameter neural network architecture, and (2) the BERT_{LARGE} model, a 24-layer, 1024-hidden, 16-heads, 340M parameter neural network architecture; both of which were trained on the [BooksCorpus](#)^[4] with 800M words, and a version of the [English Wikipedia](#) with 2,500M words.

Num of **Articles** : 5.5M

Image Grounding Ratio : **27.7%**

<Wikipedia Database>

Vokenization

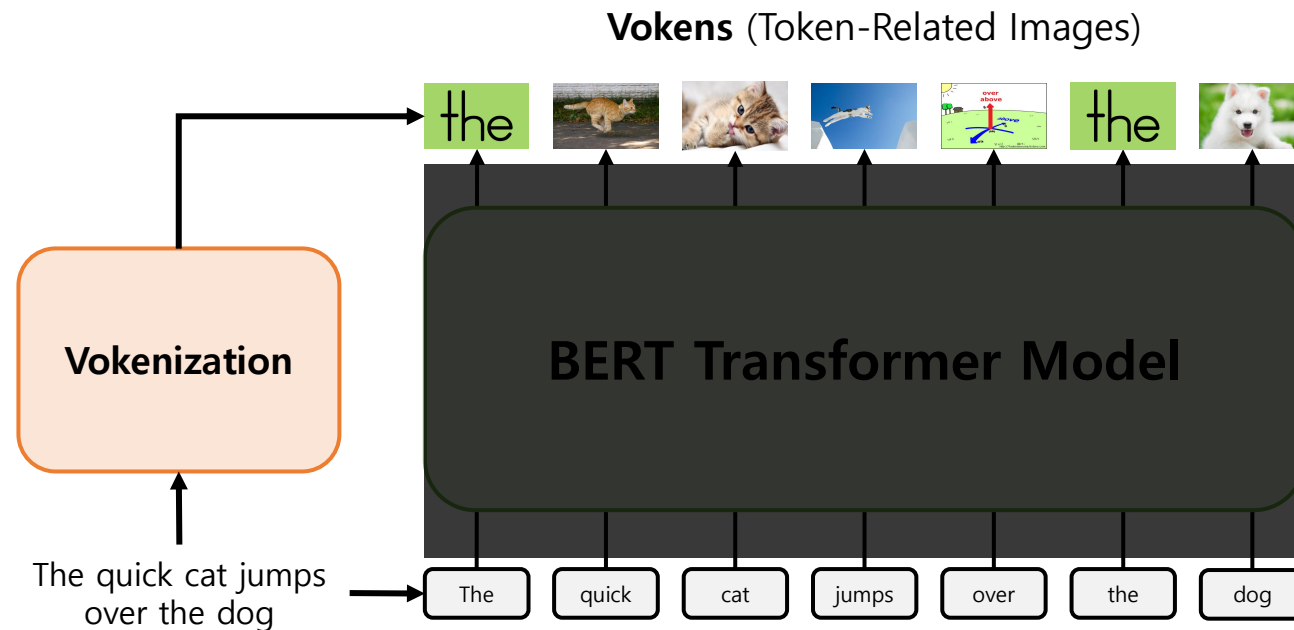
- The Vokenization Process
- Contextual Token-Image Matching Model
- Revokenization

Vokenization

-Overview

<Vokenization>

Training “**Vokenizer**” from **Image-Captioning Dataset** and
Using it to **Annotate Large Language Corpora**



Vokenization

-The Vokenization Process

<The Vokenization Process>

$w_i : \text{Token}, s = (w_1, \dots, w_l) : \text{Sentence}$

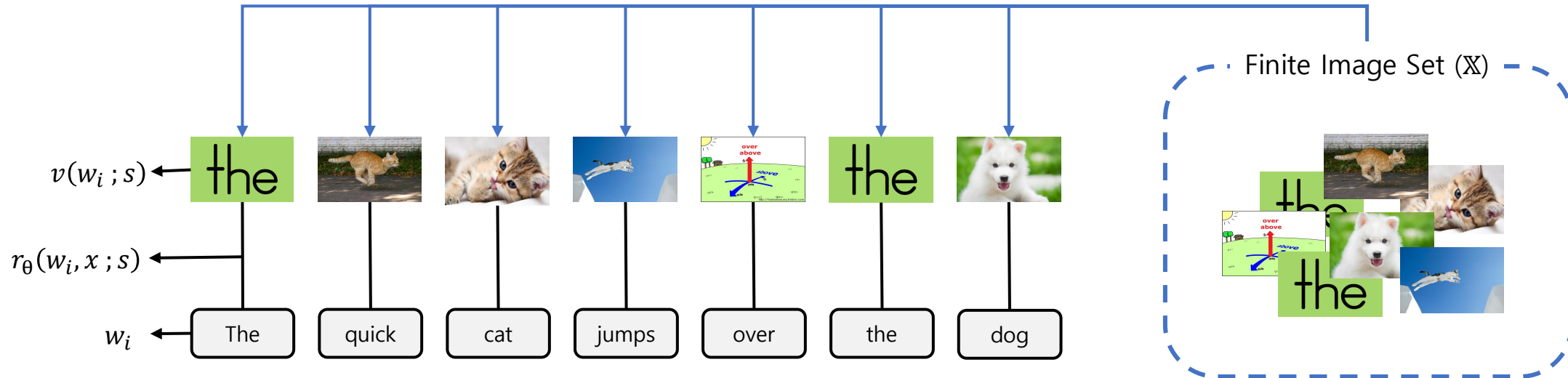
$v(w_i ; s) : \text{Image Relevant to Token (Voken)}$

$\mathbb{X} = (x_1, \dots, x_n) : \text{Finite Image Set}$

$r_\theta(w_i, x ; s) : \text{Token-Image Relevance Scoring function}$

$\theta^* : \text{Optimal Parameter}$

$$v(w_i ; s) = \operatorname{argmax}_{x \in \mathbb{X}} r_{\theta^*}(w_i, x ; s)$$



Vokenization

-The Vokenization Process

<Modeling>

$$r_{\theta}(w_i, x; s) = f_{\theta}(w_i; s)^T g_{\theta}(x) : \text{Inner Product Factorization}$$

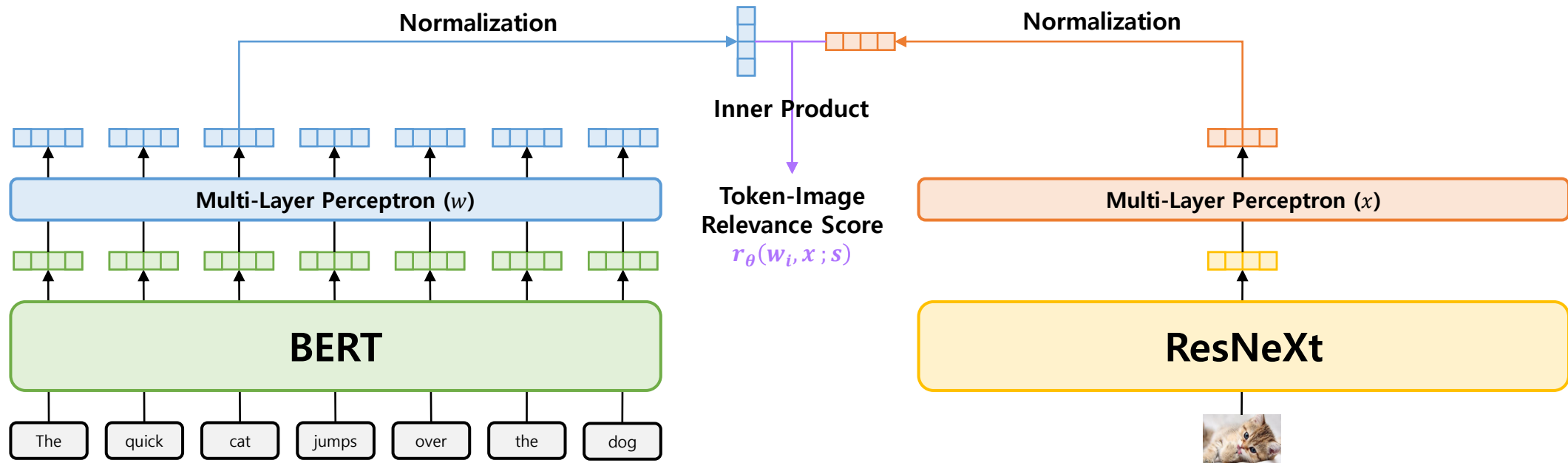
$f_{\theta}(w_i; s)$: language Feature Representation

$g_{\theta}(x)$: Image Feature Representation

$$h_1, \dots, h_l = \text{BERT}(w_1, \dots, w_l), \quad e = \text{ResNeXt}(x)$$

$w_mlp_{\theta}, x_mlp_{\theta}$: Multi-Layer Perceptron

$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{\|w_mlp_{\theta}(h_i)\|}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{\|x_mlp_{\theta}(e)\|}$$



Vokenization

-The Vokenization Process

<Modeling>

$$r_{\theta}(w_i, x; s) = f_{\theta}(w_i; s)^T g_{\theta}(x) : \text{Inner Product Factorization}$$

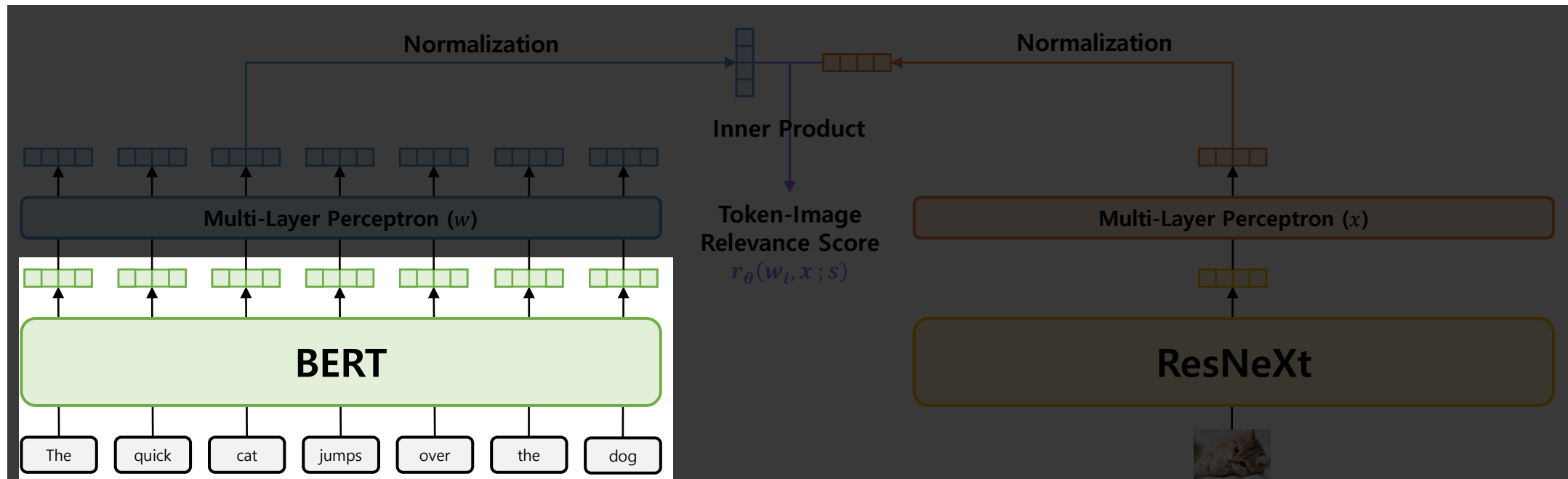
$f_{\theta}(w_i; s)$: language Feature Representation

$g_{\theta}(x)$: Image Feature Representation

$$h_1, \dots, h_l = \mathbf{BERT}(w_1, \dots, w_l), \quad e = \text{ResNeXt}(x)$$

$w_mlp_{\theta}, x_mlp_{\theta}$: Multi-Layer Perceptron

$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{||w_mlp_{\theta}(h_i)||}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{||x_mlp_{\theta}(e)||}$$



Vokenization

-The Vokenization Process

<Modeling>

$$r_{\theta}(w_i, x; s) = f_{\theta}(w_i; s)^T g_{\theta}(x) : \text{Inner Product Factorization}$$

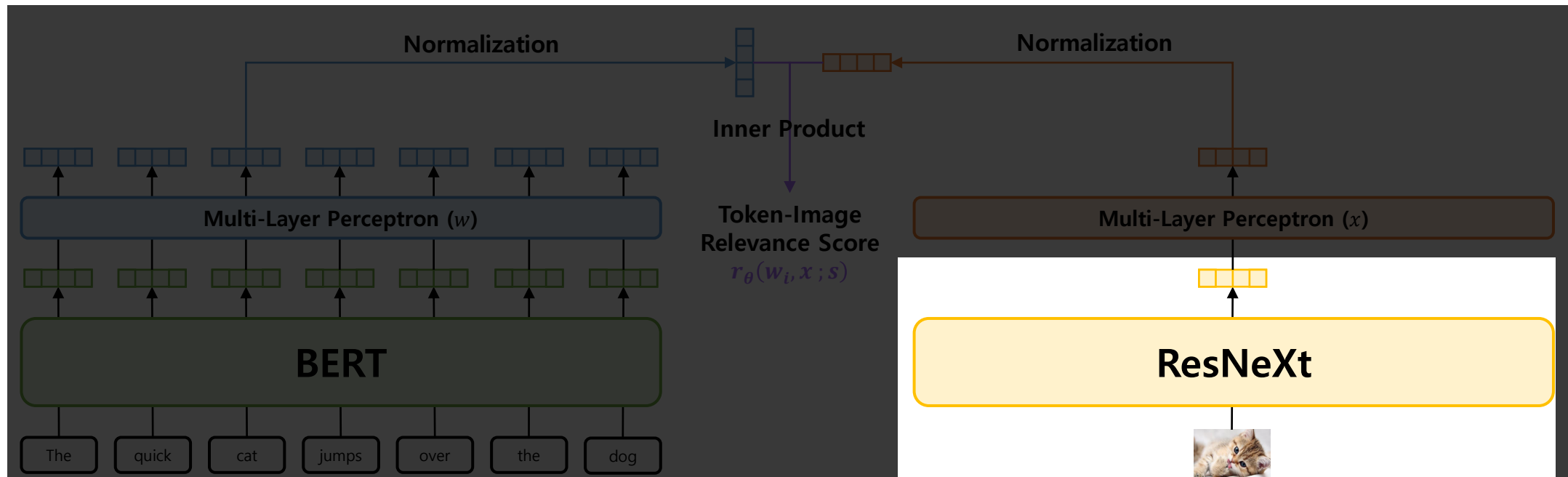
$f_{\theta}(w_i; s)$: language Feature Representation

$g_{\theta}(x)$: Image Feature Representation

$$h_1, \dots, h_l = \text{BERT}(w_1, \dots, w_l), \quad \mathbf{e} = \text{ResNeXt}(x)$$

$w_mlp_{\theta}, x_mlp_{\theta}$: Multi-Layer Perceptron

$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{\|w_mlp_{\theta}(h_i)\|}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{\|x_mlp_{\theta}(e)\|}$$



Vokenization

-The Vokenization Process

<Modeling>

$$r_{\theta}(w_i, x; s) = f_{\theta}(w_i; s)^T g_{\theta}(x) : \text{Inner Product Factorization}$$

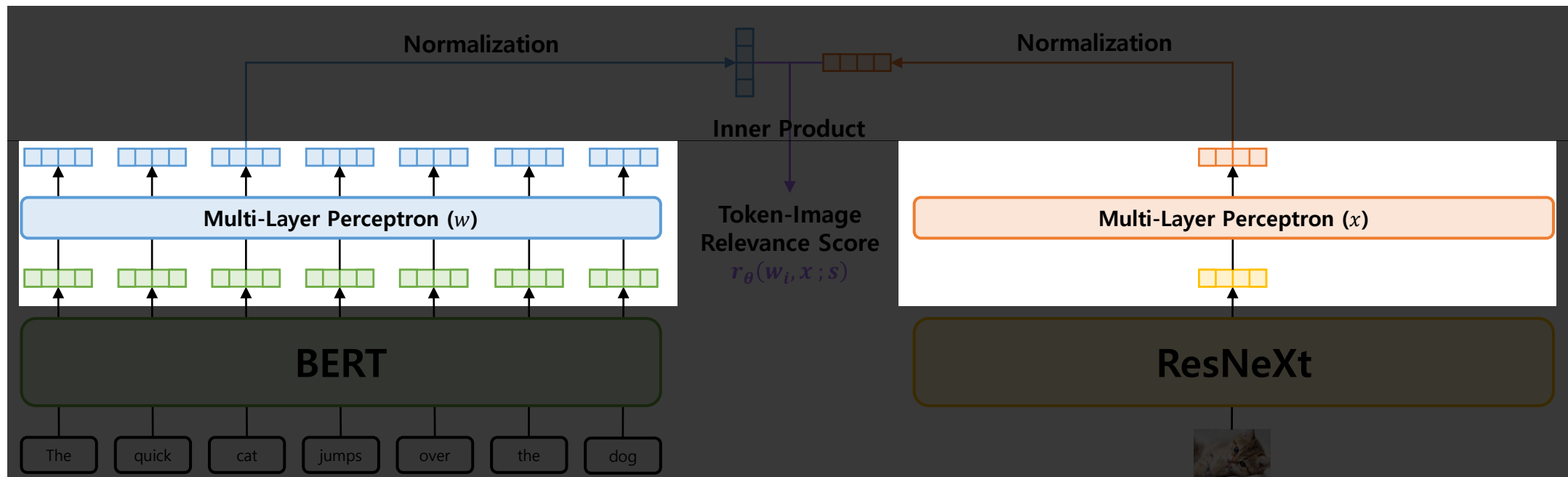
$f_{\theta}(w_i; s) : \text{language Feature Representation}$

$g_{\theta}(x) : \text{Image Feature Representation}$

$$h_1, \dots, h_l = \text{BERT}(w_1, \dots, w_l), \quad e = \text{ResNeXt}(x)$$

$w_mlp_{\theta}, x_mlp_{\theta} : \text{Multi-Layer Perceptron}$

$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{||w_mlp_{\theta}(h_i)||}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{||x_mlp_{\theta}(e)||}$$



Vokenization

-The Vokenization Process

<Modeling>

$$r_{\theta}(w_i, x; s) = f_{\theta}(w_i; s)^T g_{\theta}(x) : \text{Inner Product Factorization}$$

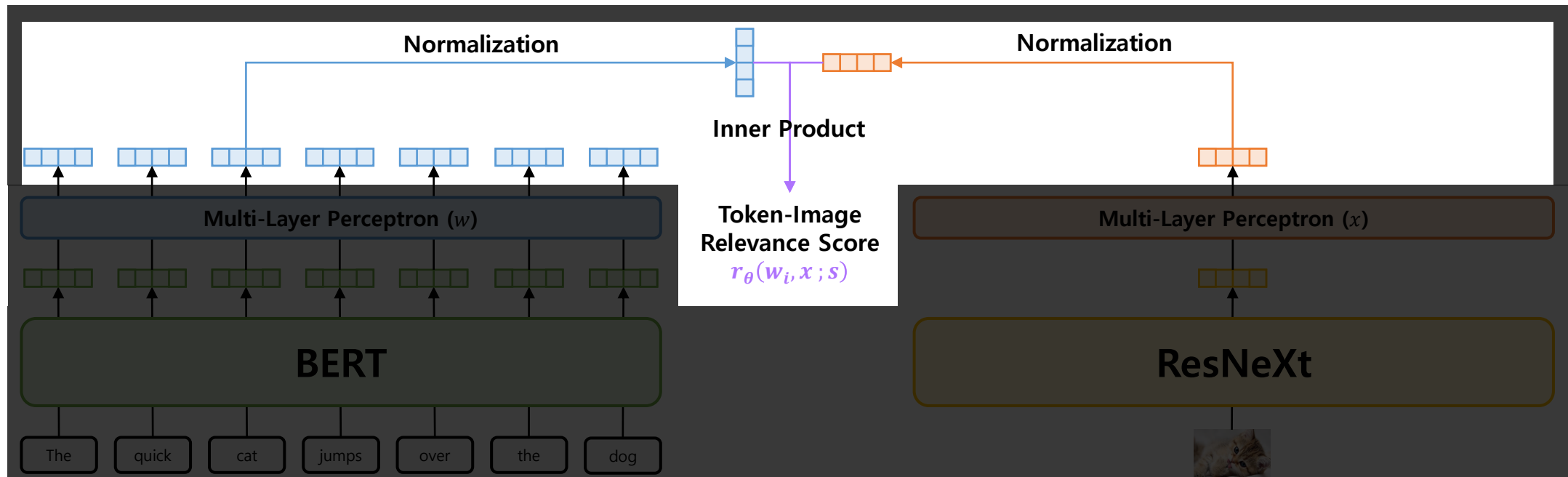
$f_{\theta}(w_i; s)$: language Feature Representation

$g_{\theta}(x)$: Image Feature Representation

$$h_1, \dots, h_l = \text{BERT}(w_1, \dots, w_l), \quad e = \text{ResNeXt}(x)$$

$w_mlp_{\theta}, x_mlp_{\theta}$: Multi-Layer Perceptron

$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{\|w_mlp_{\theta}(h_i)\|}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{\|x_mlp_{\theta}(e)\|}$$

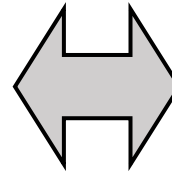


Vokenization

-The Vokenization Process

<MS COCO Dataset>

"The quick cat jumps over the dog."



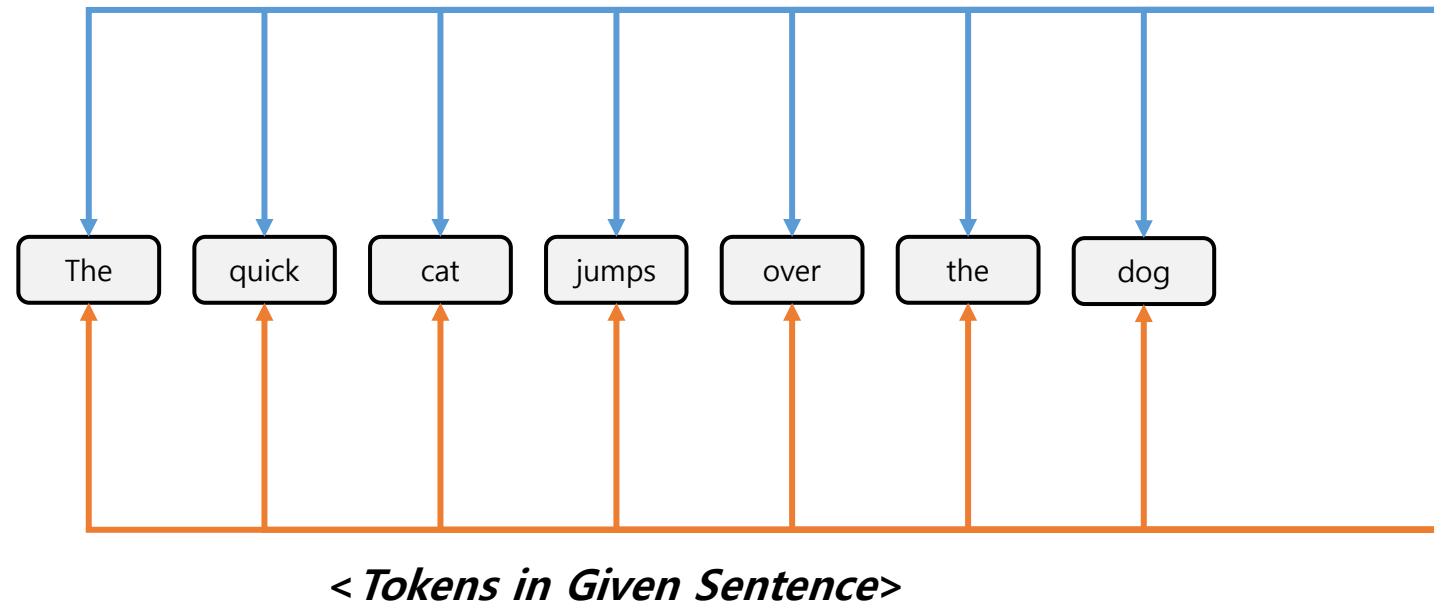
< Sentence >

< Image That Semantically Matches Given Sentence >

Vokenization

-The Vokenization Process

<Training>



< Given Image >



< Randomly Sampled Image >

Vokenization

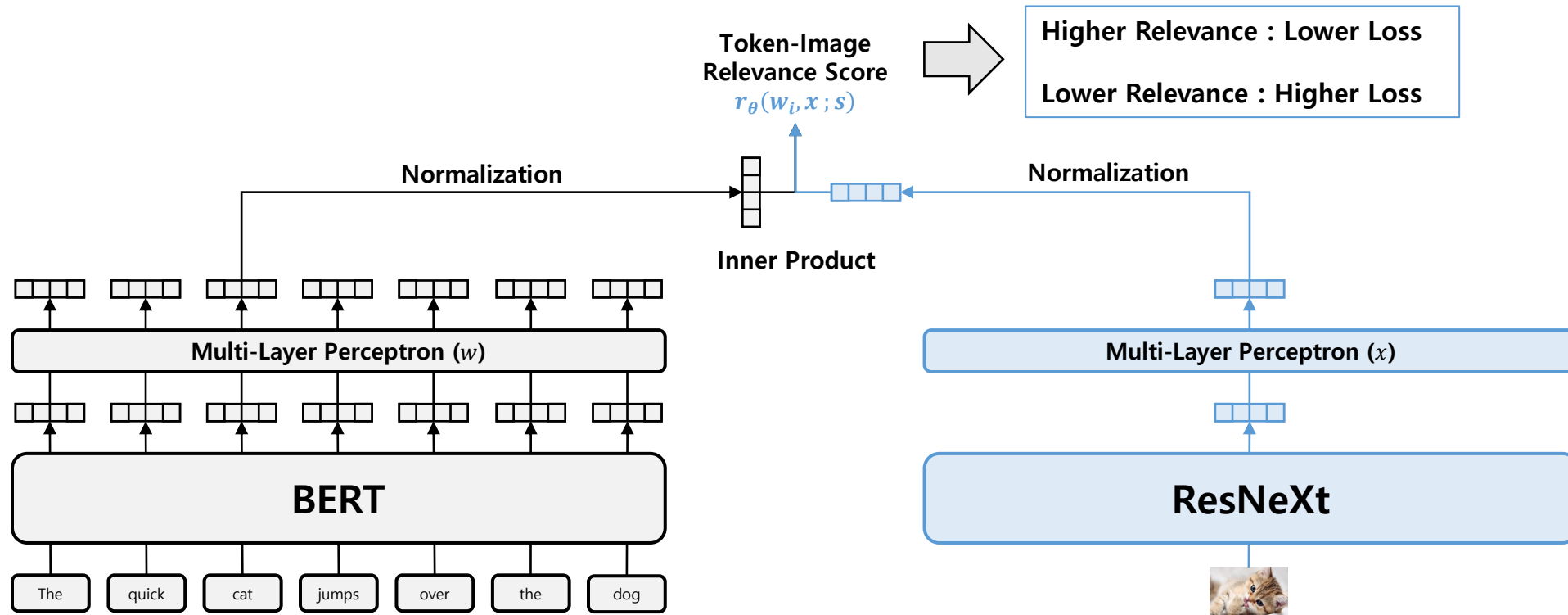
-The Vokenization Process

<Training>

(s, x) : *Ground Truth Image*–Captioning Data Point

x' : Randomly Sampled Image, where $x' \neq x$

$$\mathcal{L}_{\theta}(s, x, x') = \sum_{i=1}^l \max\{0, M - r_{\theta}(w_i, x; s) + r_{\theta}(w_i, x'; s)\}$$



Vokenization

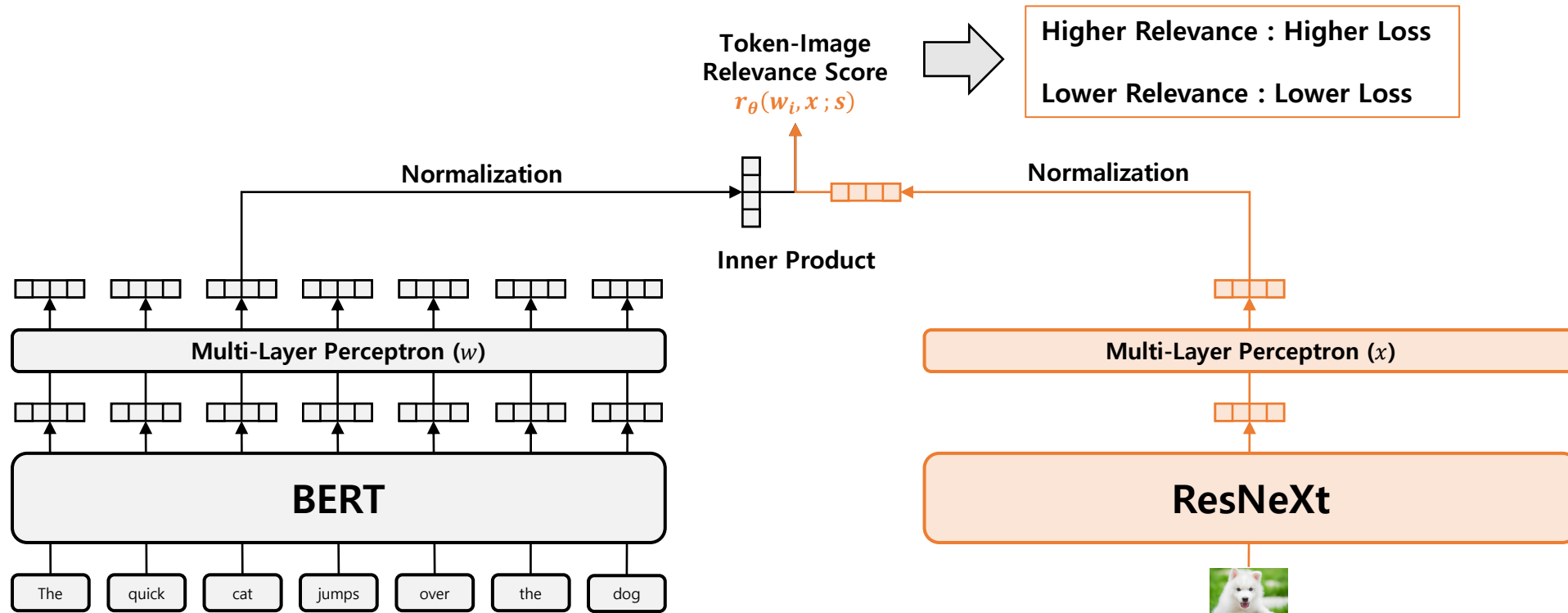
-The Vokenization Process

<Training>

(s, x) : Ground Truth Image-Captioning Data Point

x' : Randomly Sampled Image, where $x' \neq x$

$$\mathcal{L}_{\theta}(s, x, x') = \sum_{i=1}^l \max\{0, M - r_{\theta}(w_i, x; s) + r_{\theta}(w_i, x'; s)\}$$



Vokenization

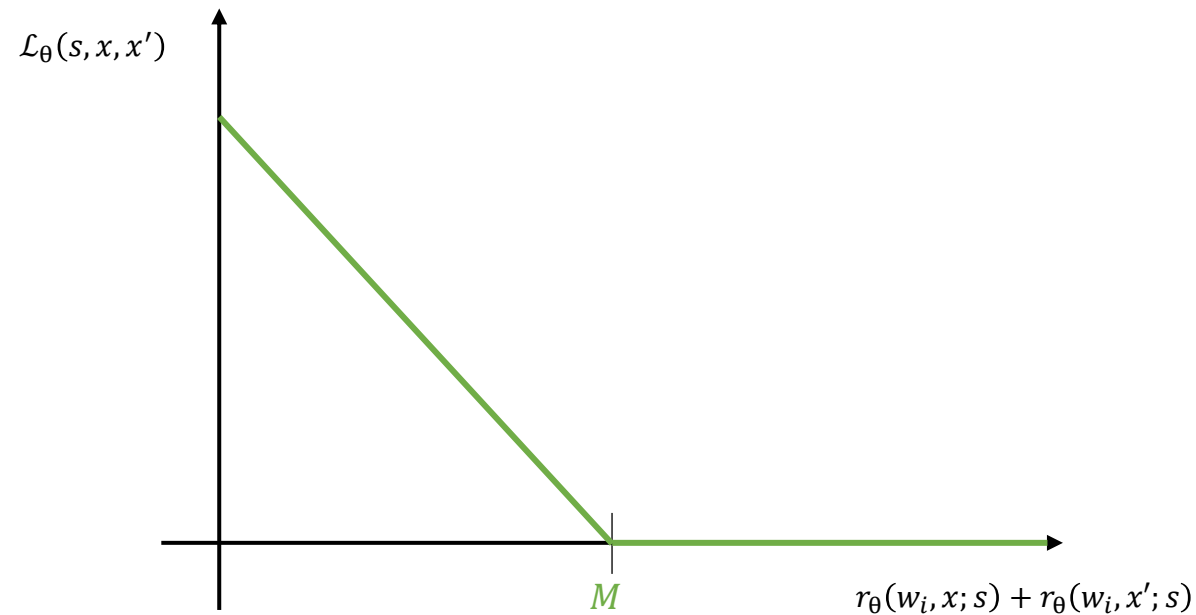
-The Vokenization Process

<Training>

(s, x) : *Ground Truth Image*–Captioning Data Point

x' : *Randomly Sampled Image*, where $x' \neq x$

$$\mathcal{L}_{\theta}(s, x, x') = \sum_{i=1}^l \max\{0, M - r_{\theta}(w_i, x; s) + r_{\theta}(w_i, x'; s)\}$$



Vokenization

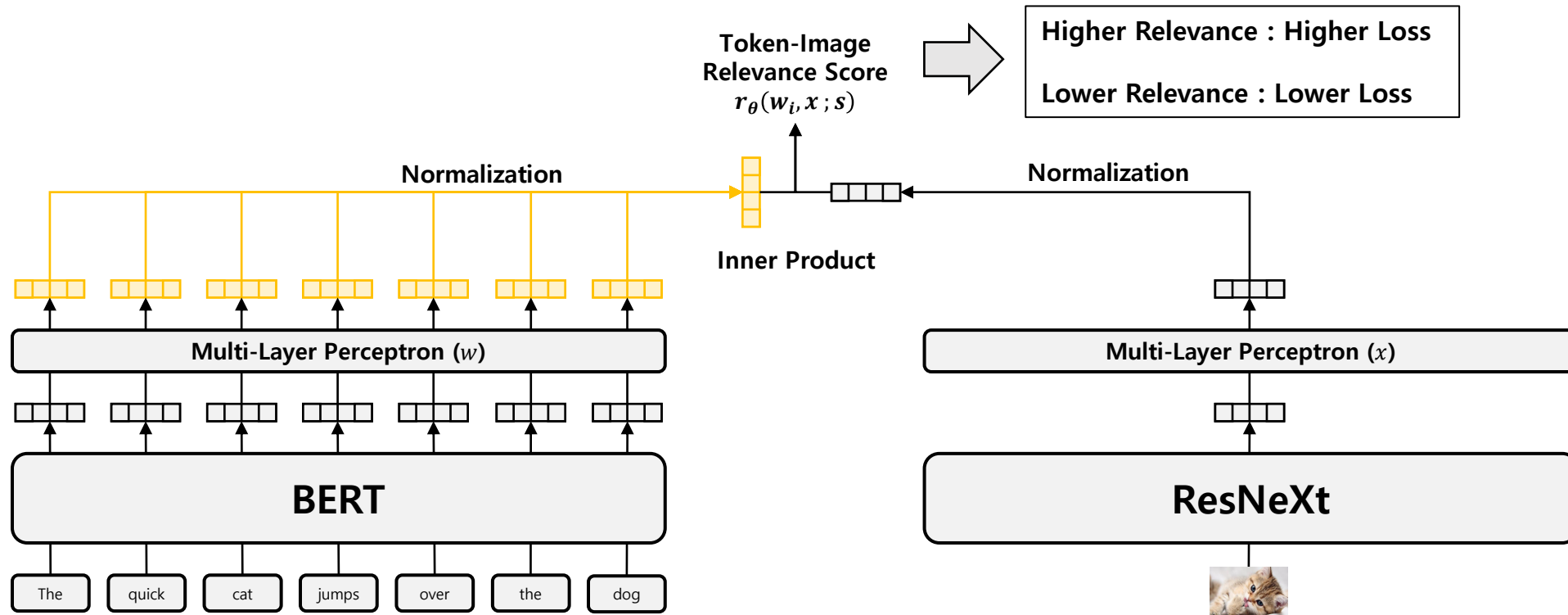
-The Vokenization Process

<Training>

(s, x) : Ground Truth Image-Captioning Data Point

x' : Randomly Sampled Image, where $x' \neq x$

$$\mathcal{L}_{\theta}(s, x, x') = \sum_{i=1}^l \max\{0, M - r_{\theta}(w_i, x; s) + r_{\theta}(w_i, x'; s)\}$$



Vokenization

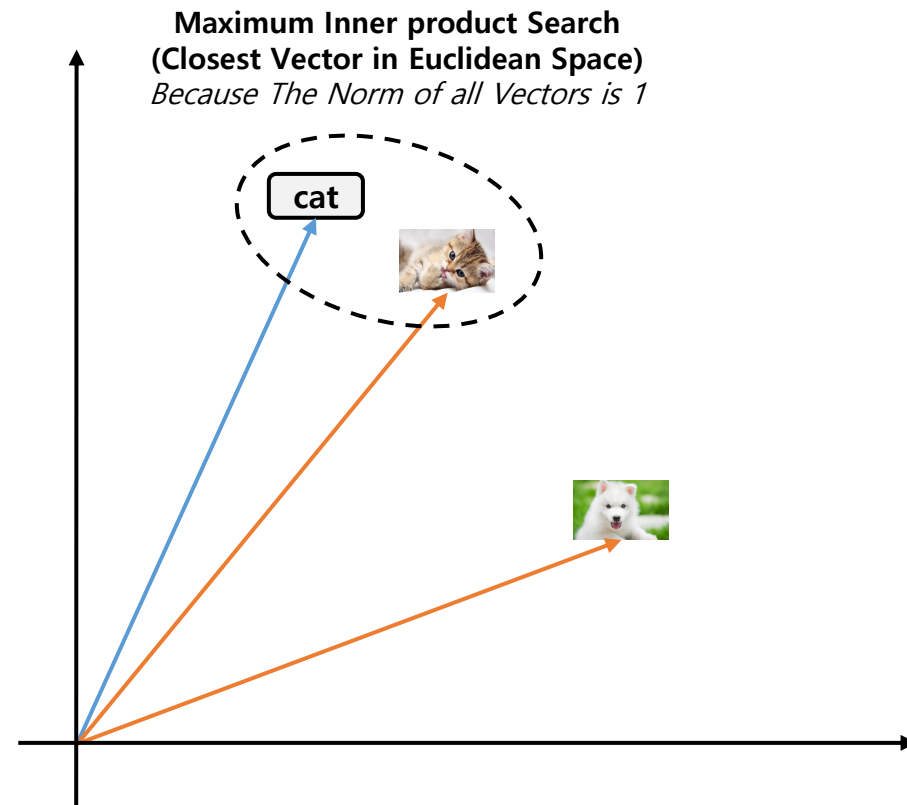
-The Vokenization Process

<Inference>

$f_{\theta}(w_i; s)$: language Feature Representation

$g_{\theta}(x)$: Image Feature Representation

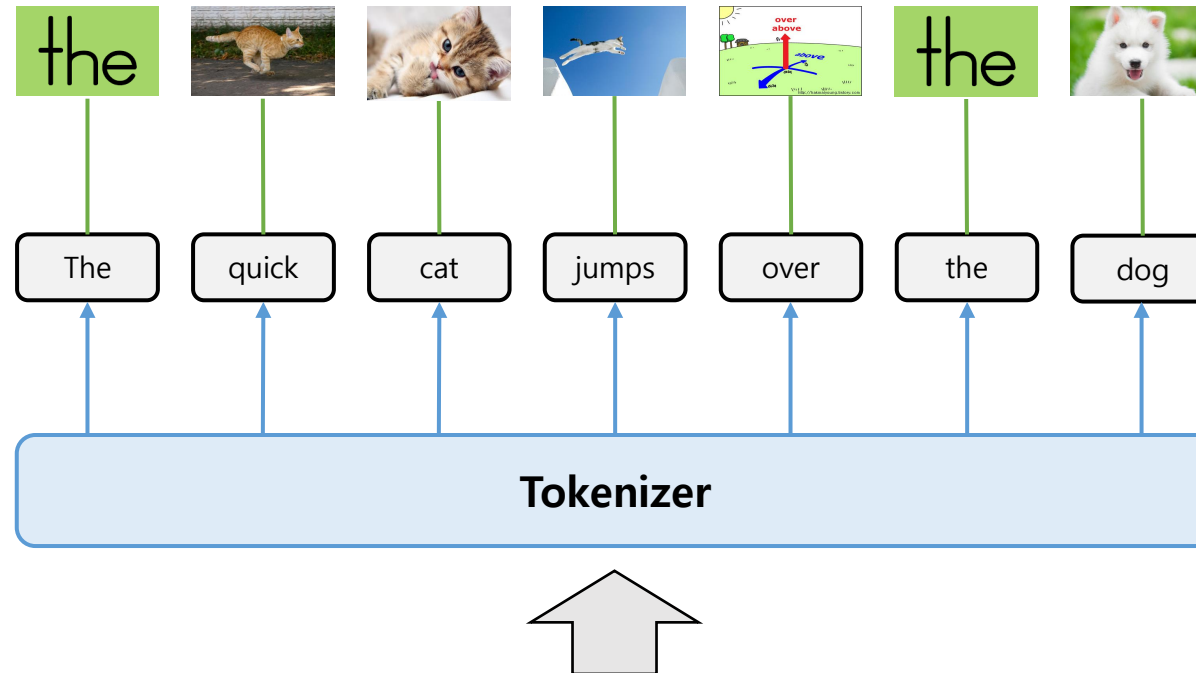
$$f_{\theta}(w_i; s) = \frac{w_mlp_{\theta}(h_i)}{\|w_mlp_{\theta}(h_i)\|}, \quad g_{\theta}(x) = \frac{x_mlp_{\theta}(e)}{\|x_mlp_{\theta}(e)\|}$$



Vokenization

-Revokenization

<Revokenization>

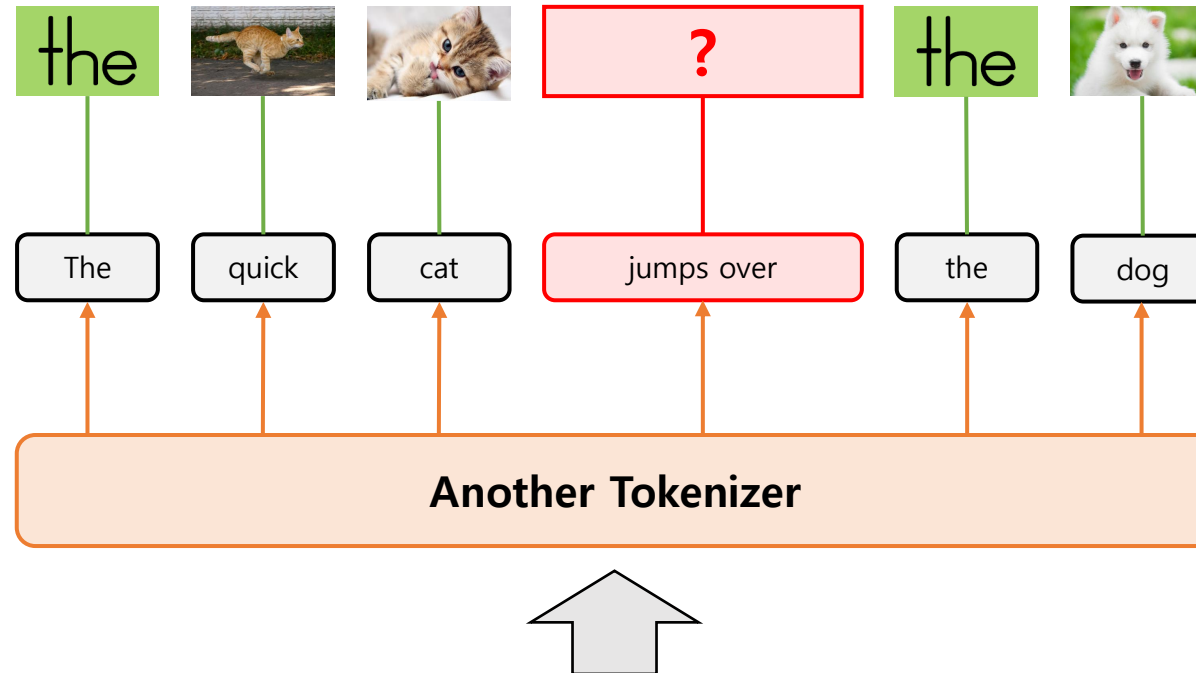


"The quick cat jumps over the dog."

Vokenization

-Revokenization

<Revokenization>



"The quick cat jumps over the dog."

<Revokenization>

$T_1, T_2 : \text{Tokenizer}$

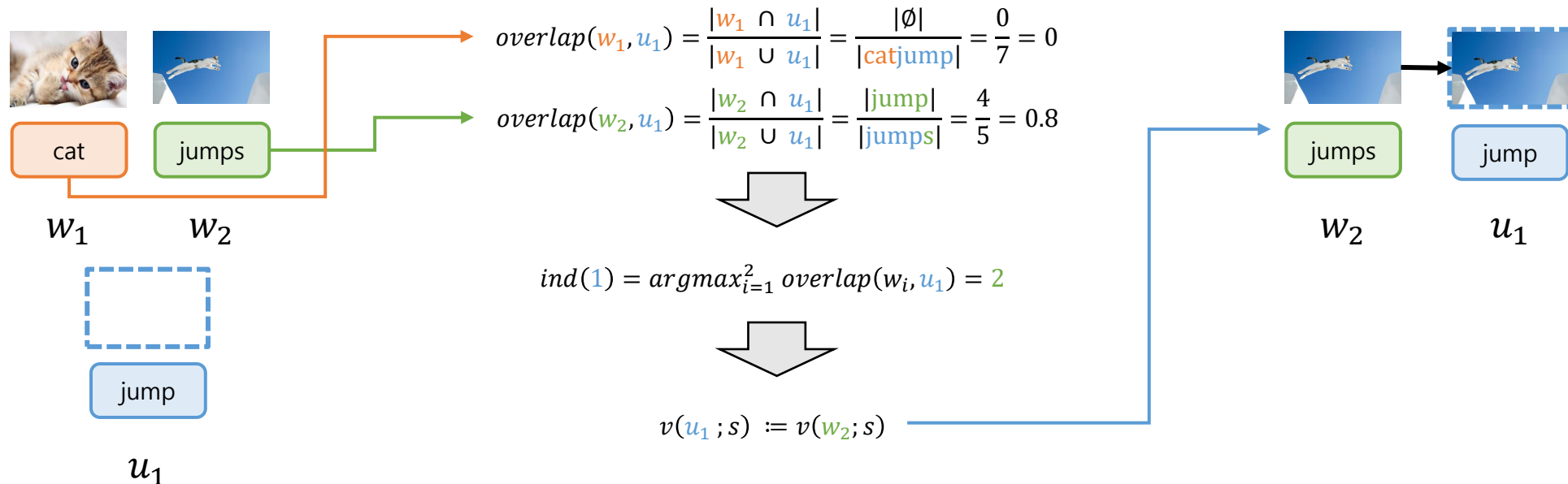
$T_1(s) = (w_1, \dots, w_l), T_2(s) = (u_1, \dots, u_m)$

$\{v(w_i; s)\}_{i=1}^l : \text{Vokens Aligned to Tokens } \{w_i\}_{i=1}^l$

$v(u_j; s) := v(w_{\text{ind}(j)}; s)$

$\text{ind}(j) = \text{argmax}_{i=1}^l \text{overlap}(w_i, u_j)$

$\text{overlap}(w_i, u_j) = \frac{|w_i \cap u_j|}{|w_i \cup u_j|}$



Experiments

- Data, Task & Details
- Result

Experiments

-Data, Task & Details

<Data, Task & Details>

<Pre-Training Data>

English Wikipedia, Wiki103

<Fine-Tuning Task>

GLUE, SQuAD, SWAG, SST-2, QNLI, QQP, MNLI

<Implementation Details>

Concatenation of last 4 layers of BERT

ResNeXt-101-32x8d

Two Fully-Connected Multi-Layer Perceptron, 256-dim ($w_{mlp_\theta}, x_{mlp_\theta}$)

ReLU, 64-dim output

$M = 0.5$, $\lambda = 1$, vocab_size = 50,000, batch_size = 256, epochs = 3

Removing Next-Sentence Prediction

12-layer BERT_{BASE}, 768-dim, **From Scratch**

Experiments

-Result

<Result>

Method	STT-2	QNLI	QQP	MNLI	SQuAD v1.1	SQuAD v2.0	SWAG	Avg.
BERT _{6L/512H}	88.0	85.2	87.1	77.9	71.3/80.2	57.2/60.8	56.2	75.6
BERT _{6L/512H} + Voken-cls	89.7	85.0	87.3	78.6	71.5/80.2	61.3/64.6	58.2	76.8
BERT _{12L/768H}	89.3	87.9	83.2	79.4	77.0/85.3	67.7/71.1	65.7	79.4
BERT _{12L/768H} + Voken_cls	92.2	88.6	88.6	82.6	78.8/86.7	68.1/71.2	70.6	82.1
RoBERTa _{6L/512H}	87.8	82.4	85.2	73.1	50.9/61.9	49.6/52.7	55.1	70.2
RoBERTa _{6L/512H} + Voken_cls	87.8	85.1	85.3	76.5	55.0/66.4	50.9/54.1	60.0	72.6
RoBERTa _{12L/768H}	89.2	87.5	86.2	79.0	70.2/79.9	59.2/63.1	65.2	77.6
RoBERTa _{12L/768H} + Voken_cls	90.5	89.2	87.8	81.0	73.0/82.5	65.9/69.3	70.4	80.6

Experiments

-Result

<Result>

Method	STT-2	QNLI	QQP	MNLI	SQuAD v1.1	SQuAD v2.0	SWAG	Avg.
BERT _{6L/512H}	88.0	85.2	87.1	77.9	71.3/80.2	57.2/60.8	56.2	75.6
BERT _{6L/512H} + Voken-cl	89.7	85.0	87.3	78.6	71.5/80.2	61.3/64.6	58.2	76.8
BERT _{12L/768H}	89.3	87.9	83.2	79.4	77.0/85.3	67.7/71.1	65.7	79.4
BERT _{12L/768H} + Voken_cls	92.2	88.6	88.6	82.6	78.8/86.7	68.1/71.2	70.6	82.1
RoBERTa _{6L/512H}	87.8	82.4	85.2	73.1	50.9/61.9	49.6/52.7	55.1	70.2
RoBERTa _{6L/512H} + Voken_cls	87.8	85.1	85.3	76.5	55.0/66.4	50.9/54.1	60.0	72.6
RoBERTa _{12L/768H}	89.2	87.5	86.2	79.0	70.2/79.9	59.2/63.1	65.2	77.6
RoBERTa _{12L/768H} + Voken_cls	90.5	89.2	87.8	81.0	73.0/82.5	65.9/69.3	70.4	80.6

Analysis

- Visualization of Vokens

Analysis

-Visualization of Vokens

<Visualization of Vokens>



"Humans learn language by listening, speaking, writing, reading"



"Down by the salley gardens my love and I did meet"

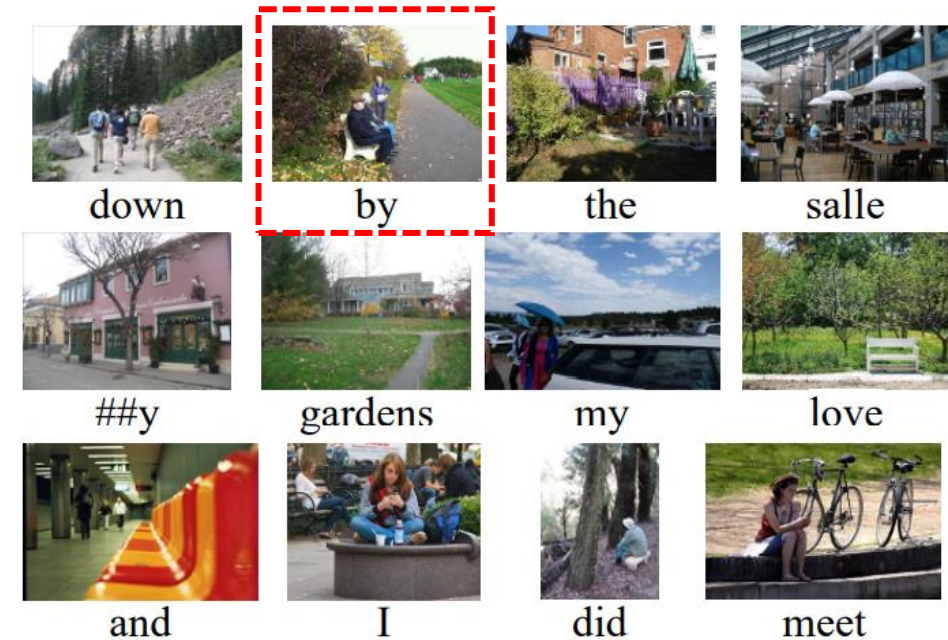
Analysis

-Visualization of Vokens

<Visualization of Vokens>



"Humans learn language by listening, speaking, writing, reading"



"Down by the salley gardens my lobe and I did meet"

*The Vokenizer Maps the Voken in Consideration of the Context
and Can Map Vokens to Non-Concrete Tokens*

Analysis

-Visualization of Vokens

<Visualization of Vokens>



"Humans learn language by listening, speaking, writing, reading"



"Down by the salley gardens my lobe and I did meet"

Some Vokens Has a Shift in Alignment Which May be Caused by Limitations of Sentence-Image Weak Supervision

Analysis

-Visualization of Vokens

<Visualization of Vokens>



"Humans learn language by listening, speaking, writing, reading"



"Down by the salley gardens my lobe and I did meet"

The Related Visual Information Can Help Understand the Language and Lead to the Improvement

Conclusion

<Conclusion>

- **This Paper Explored the Possibility of Utilizing Visual Supervision to Language Encoder.**
- **Vokenizer is Contextual Token-Image Matching Model and Can be Used to Vokenize the Language Corpus.**
- **Supervised by These Generated Vokens, The Language Model Has a Significant Improvement Over the Purely Self-Supervised Language Model on Multiple Language Tasks.**

Any Questions?

Thank You