# Seq2seq & Attention

## 2023-2 KUBIG 방학세션
## DL

# 0. Announcement

## No real-time Class on Next Week!

$\Rightarrow$ Instead, Pytorch implementation of Transformer

- I would give you the youtube link on next week.
- Full course of transformer
- Refered to "Attention is All you need"
- GPU?
- After next week, we would prepare "KUBIG Contest"

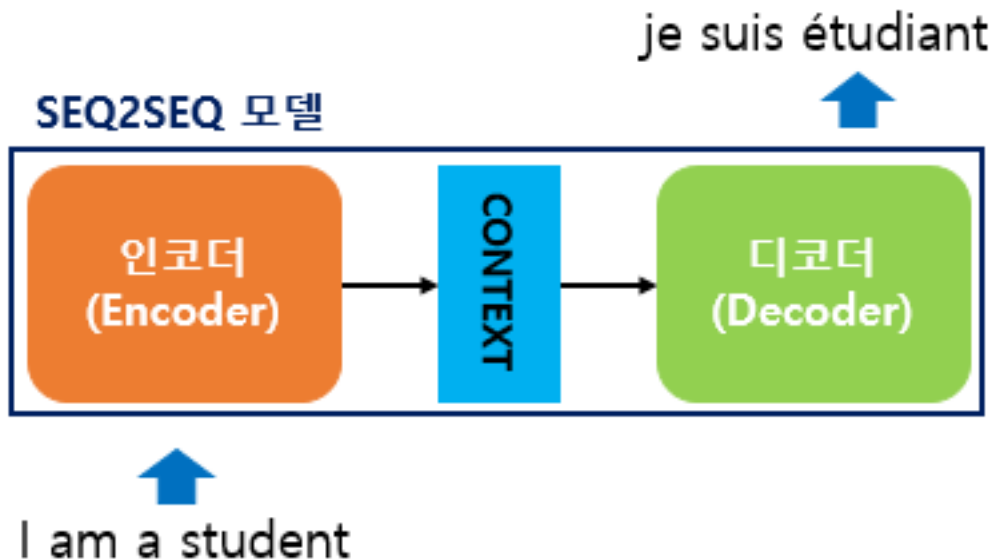https://dacon.io/competitions/official/235747/overview/description

# Category

1. Seq2Seq
2. Attention
3. Training Techniques
4. Q&A

# 1. Seq2Seq

# 1. Seq2Seq

## Many-to-Many: Machine Translation | Chatbot

je suis étudiant

SEQ2SEQ 모델

인코더
(Encoder)

CONTEXT

디코더
(Decoder)

I am a student

# 1. Seq2Seq

## Many-to-Many: Machine Translation | Chatbot

# 1. Seq2Seq

## Word Embedding: More powerful Strategy

| I | | | | am | | | | a | | | | student | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.157 | | | | 0.78 | | | | 0.75 | | | | 0.88 |
| | | -0.25 | | | | 0.29 | | | | -0.81 | | | | -0.17 |
| | | 0.478 | | | | -0.96 | | | | 0.96 | | | | 0.29 |
| | | -0.78 | | | | 0.52 | | | | 0.12 | | | | 0.48 |

# 1. Seq2Seq

1) Tokenize

"KUBIG is very good!"
-> "KUBIG" "is" "very" "good" "!"

2) Punctuation

-> "KUBIG" "is" "very" "good"

-> Removing is always not a good method!

: It's not enough to make word embeddings!

# 1. Seq2Seq

## 1) Tokenize

```python
from nltk.tokenize import word_tokenize
from nltk.tokenize import WordPunctTokenizer
```

```
['Don', "'", 't', 'be', 'fooled', 'by', 'the', 'dark', 'sounding', 'name', ',', 'Mr', '.', 'Jone',
"'", 's', 'Orphanage', 'is', 'as', 'cheery', 'as', 'cheery', 'goes', 'for', 'a', 'pastry', 'shop',
'.']
```

## 2) Punctuation

```python
from nltk.tokenize import RegexpTokenizer
```

```
['Think', 'and', 'analyze', 'analyze', 'and', 'think', 'The', 'loop']
```

# 1. Seq2Seq

How about this situation?

"Don't use that spoon"

1) Do / n't / use / that / spoon

2) Don' / t / use / that / spoon

# 1. Seq2Seq

**How about this situation?**

"Don't use that spoon"

1) Do / n't / use / that / spoon  -> Ordinary Usage

2) Don' / t / use / that / spoon

# 1. Seq2Seq

## Word Preprocessing

1) Case conversion
   -> Make all words into small cases

2) Cleansing noise data
   - infrequent words
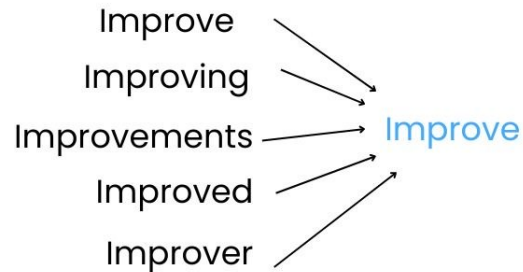   - short words('a', 'an', 'is' etc)

# 1. Seq2Seq

## Word Preprocessing
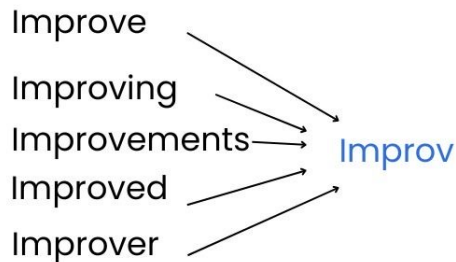
### 1) Stemming
  : remove few characters

### 2) Lemmatization
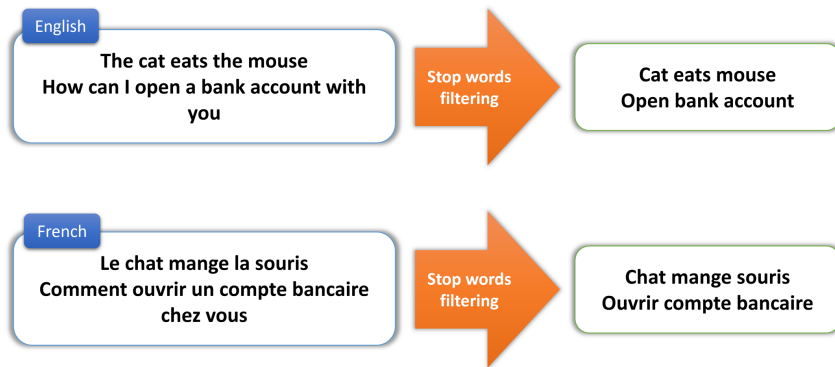  : conversion into meaningful base

**Stemming vs Lemmatization**

Improve
Improving
Improvements → Improv
Improved
Improver

Improve
Improving
Improvements → Improve
Improved
Improver

# 1. Seq2Seq

## Word Preprocessing

### 3) Stopwords

: Actually, 'me, I, us' are not meaningful when we analyze the real natural languages

: Other unique words are more attentive!

English
| The cat eats the mouse
How can I open a bank account with you |

Stop words filtering →

Cat eats mouse
Open bank account

French
| Le chat mange la souris
Comment ouvrir un compte bancaire chez vous |

Stop words filtering →

Chat mange souris
Ouvrir compte bancaire

# 1. Seq2Seq

## Word Preprocessing

### 1) Stemming

```python
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize


stemmer = PorterStemmer()
```

### 2) Lemmatization

```python
from nltk.stem import WordNetLemmatizer


lemmatizer = WordNetLemmatizer()
```

### 3) Stopwords

```python
from nltk.corpus import stopwords
```

NLTK

# 1. Seq2Seq

## Word Preprocessing

## => Integer Encoding

: Computer couldn't understand our real language -> We have to convert them into number!

1) Index Encoding
2) One-hot Encoding

The quick brown fox jumped over the brown dog

| the | quick | brown | fox | jumped | over | the | brown | dog |
|-----|-------|-------|-----|--------|------|-----|-------|-----|
| 1 | 4 | 13 | 9 | 5 | 2 | 1 | 13 | 23 |

Number of words in document

# 1. Seq2Seq

## Word Preprocessing

## => Integer Encodir

: Computer couldn't understar                                    hem into number!

1) Index Encoding
2) One-hot Encoding

-> Simple Python function
   could make Embeddings!

One-Hot Encoding

The quick brown fox jumped over the brown dog

| | cat | the | quick | brown | fox | jumped | over | dog | bird | flew | ... | kangaroo | house |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| time | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | ... | 0 | 0 |

Dictionary Size

# 1. Seq2Seq

## Word Preprocessing

# +) Padding

: Fitting the dimension is required for parallel operation

: zero-pad for short sentences

```python
for sentence in encoded:
    while len(sentence) < max_len:
        sentence.append(0)

padded_np = np.array(encoded)
padded_np
```
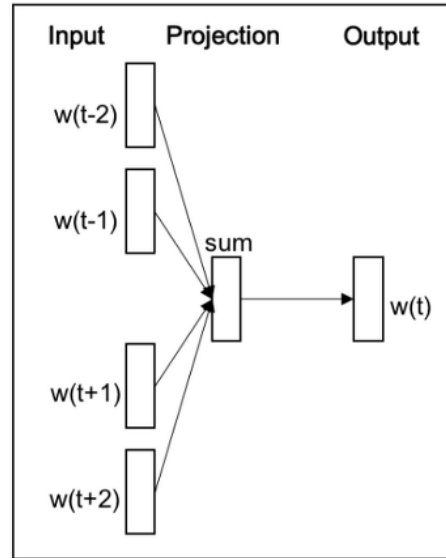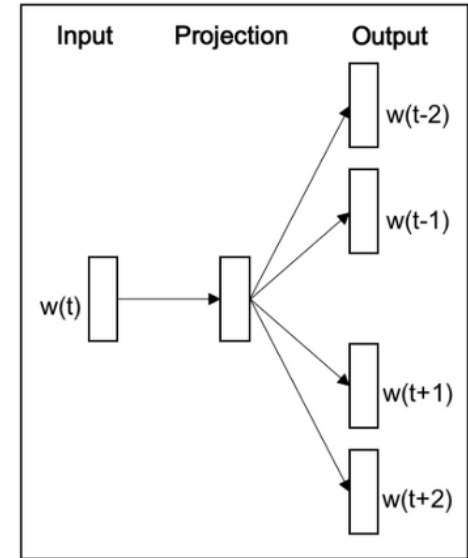
# 1. Seq2Seq

## Word Preprocessing

## +) Word2Vec

: Library for making Natural words to Vectors

: Training is necessary but powerful

- Feature
  : Distributed Expression ←→ sparse
  : Skip-gram
  : CBOW

Word2vec

king ⟶ 

man ⟶

woman ⟶

# 1. Seq2Seq

## Word Preprocessing

## +) Word2Vec

: Library for making Natural words to Vec

: Training is necessary but powerful

- Feature
  : Distributed Expression ←→ sparse
  : Skip-gram
  : CBOW



**CBOW**

| Input | Projection | Output |
|-------|-----------|--------|
| w(t-2) | | |
| w(t-1) | sum | w(t) |
| w(t+1) | | |
| w(t+2) | | |

**Skip-Gram**

| Input | Projection | Output |
|-------|-----------|--------|
| | | w(t-2) |
| w(t) | | w(t-1) |
| | | w(t+1) |
| | | w(t+2) |

# 1. Seq2Seq

```
from gensim.models import Word2Vec
model = Word2Vec(split, size=100, window=5, min_count=20, workers=4,
iter=50, sg=1)
```

# Word Preprocessing

# +) Word2Vec

: Library for making Natural words to Vec

: Training is necessary but powerful

- Feature
  : Distributed Expression ←→ sparse
  : Skip-gram
  : CBOW



**CBOW**

Input   Projection   Output

w(t-2)

w(t-1)

sum

w(t)

w(t+1)

w(t+2)

**Skip-Gram**

Input   Projection   Output

w(t)

w(t-2)

w(t-1)

w(t+1)

w(t+2)

# 1. Seq2Seq

## Seq2Seq Limitation

- Trying to compress and summarize all the information in the input sequence into one fixed-sized vector (context vector), inevitably leads to loss of information

- Gradient vanishing/exposing phenomena inevitably occur

Sigmoid function and it's derivative:

$\sigma(x)$
$\frac{d}{dx}\sigma(x)$

Saturating

Saturating

Derivative
Close to 0

Derivative
Close to 0

# 2. Attention

## Pre-Assumption

-> The decoder hidden state immediately before the decoder outputs the word will be like the encoder hidden state immediately after the encoder reads the word deeply related to in the input sequence
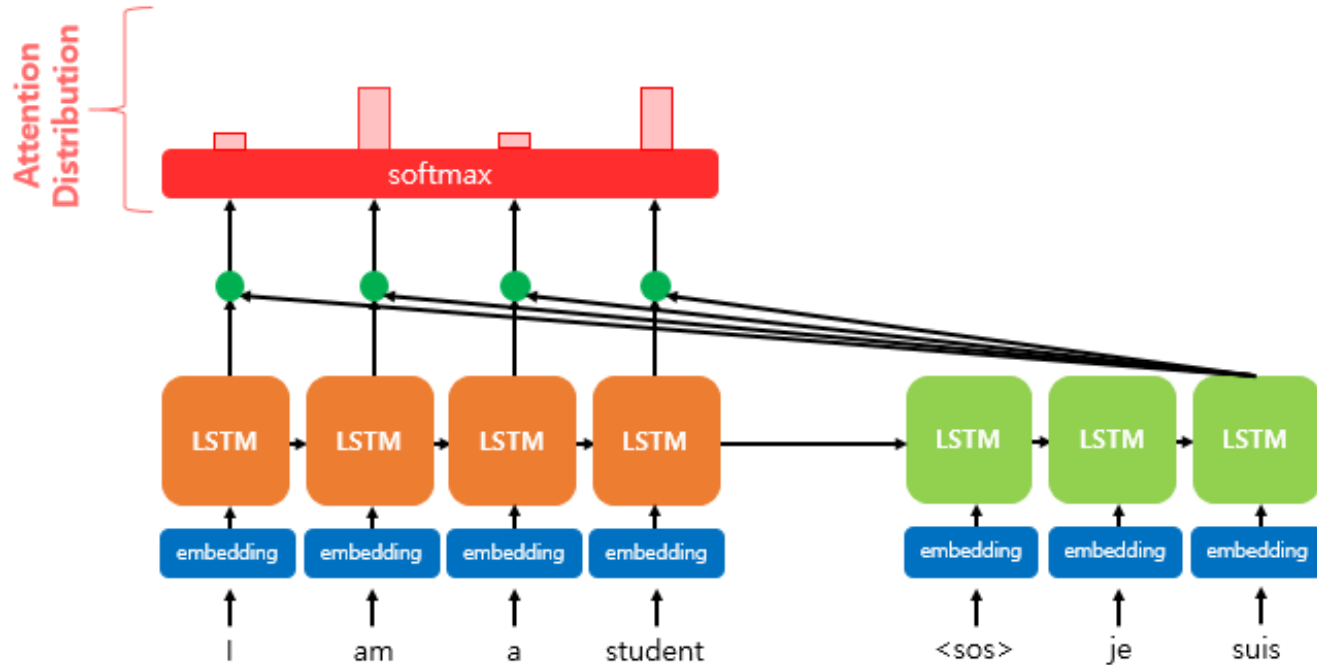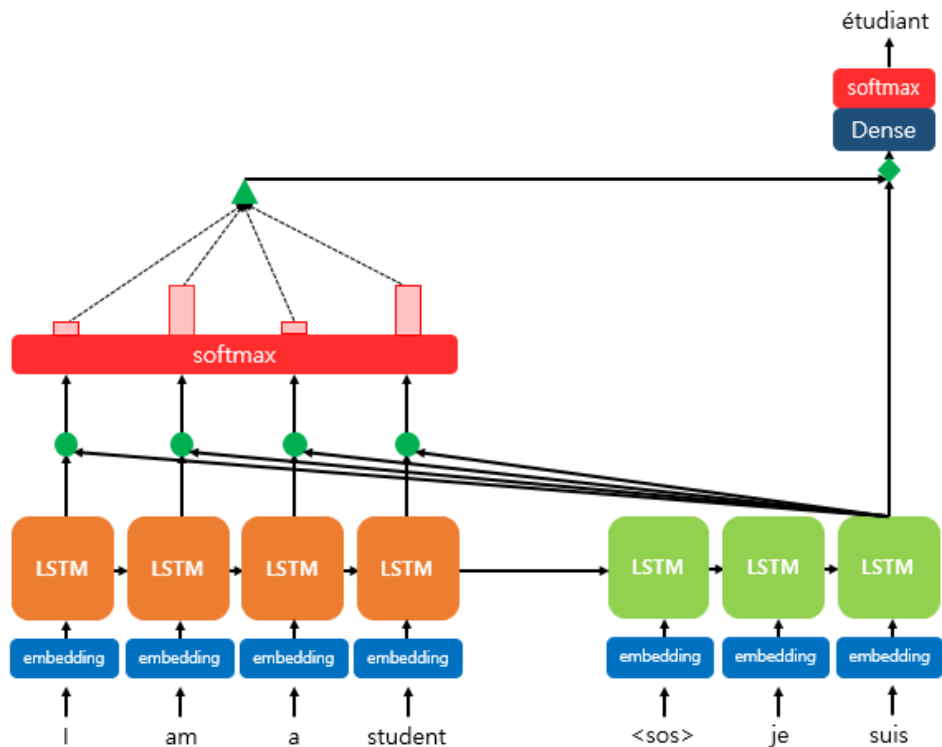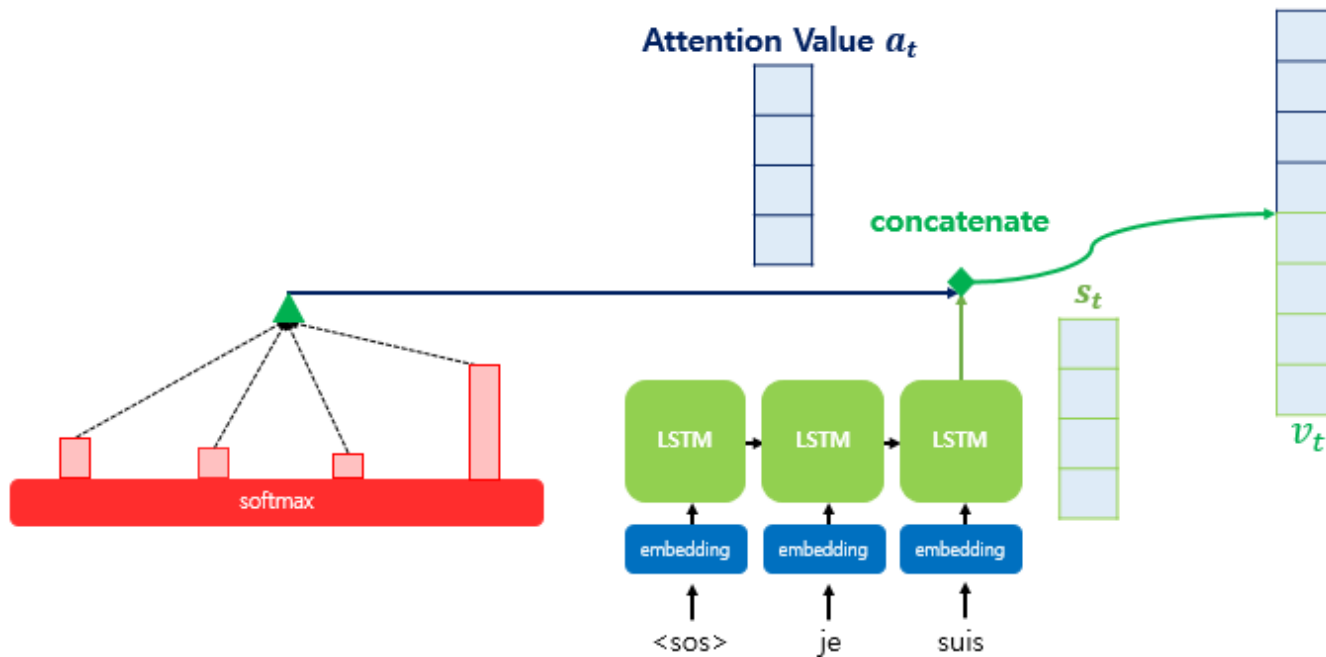
*Use all hidden state onto each decoder layer!

# 2. Attention

## Attention Mechanism

# 2. Attention

## Attention Mechanism

# 2. Attention

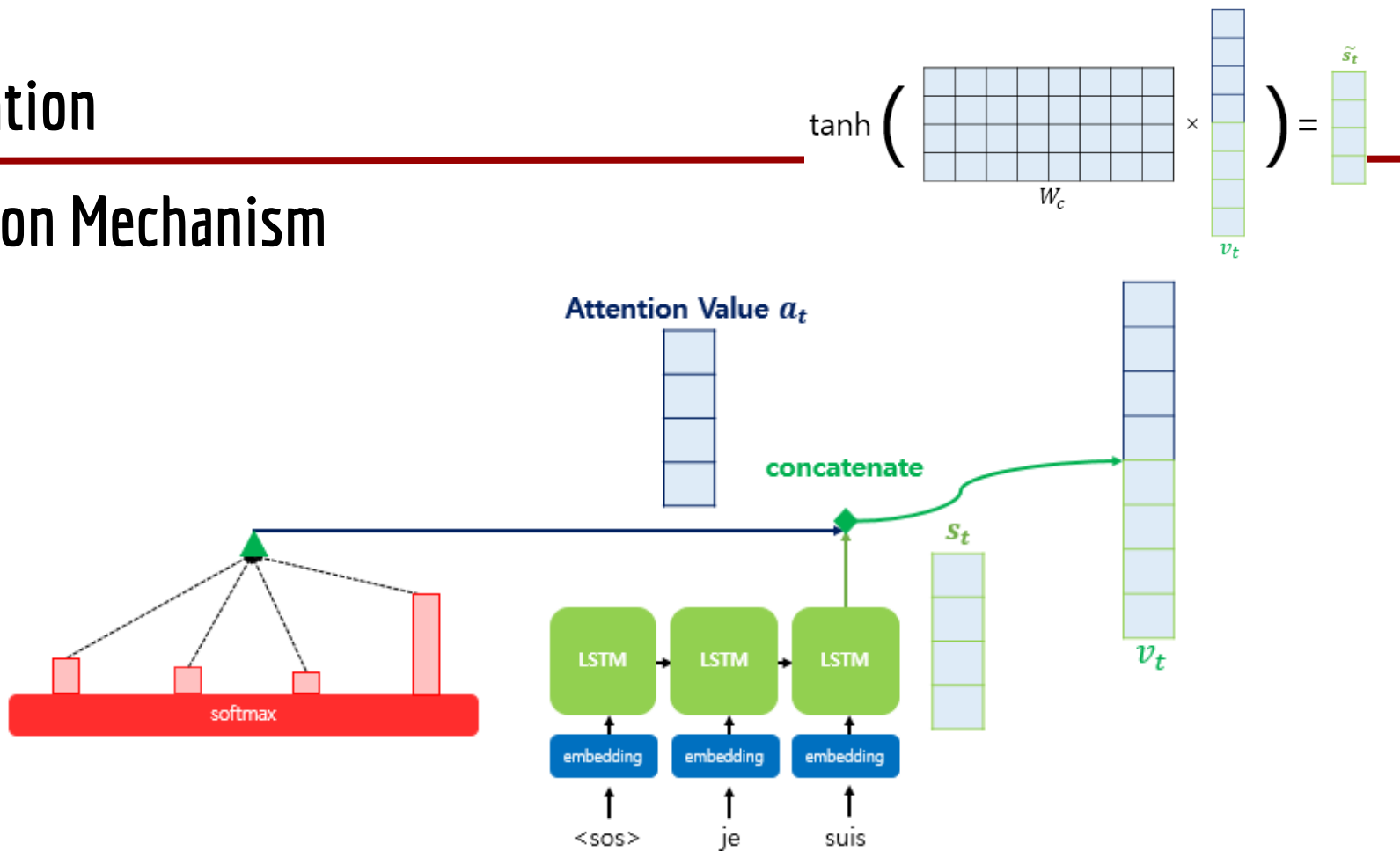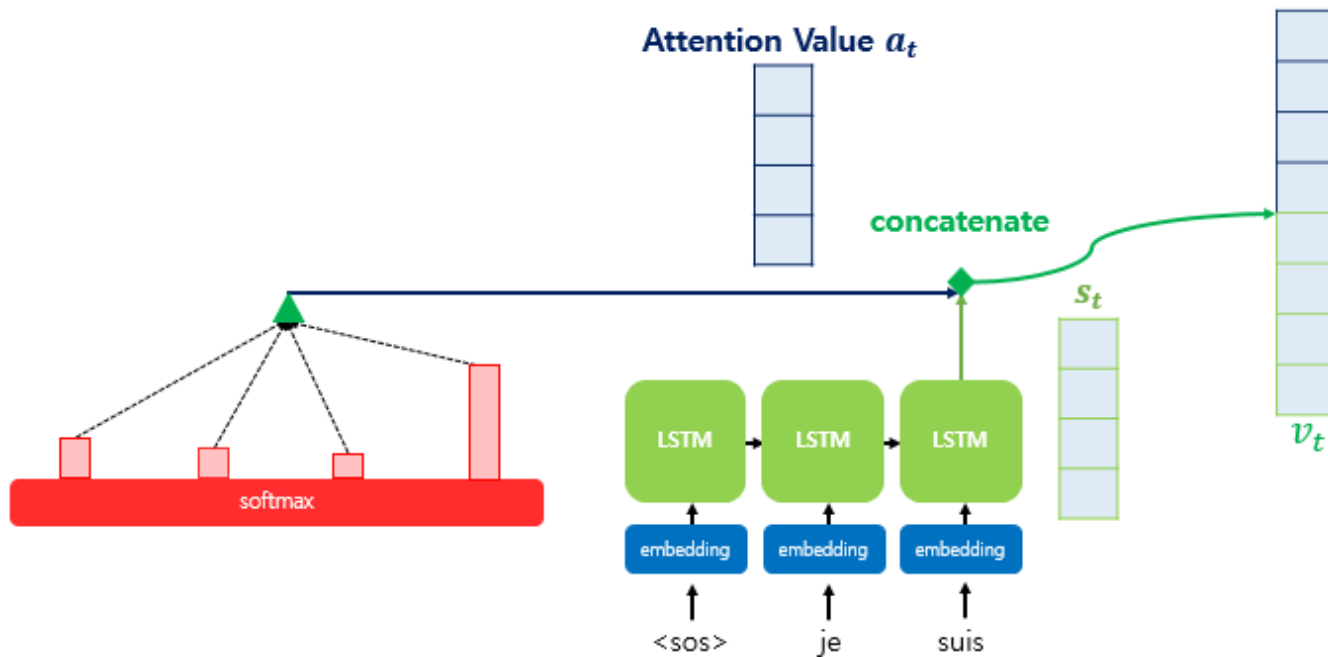## Attention Mechanism

# 2. Attention

## Attention Mechanism



Physics of dot product

# 2. Attention

## Attention Mechanism



$$score(s_t, \; h_i) = s_t^T h_i$$

# 2. Attention

## Attention Mechanism

# 2. Attention

## Attention Mechanism

# 2. Attention

## Attention Mechanism

# 2. Attention

## Attention Mechanism

$$\tanh \left( \boxed{W_c} \times \boxed{v_t} \right) = \boxed{\tilde{s}_t}$$

**Attention Value** $a_t$

**concatenate**

$s_t$

$v_t$

softmax

| LSTM | LSTM | LSTM |

| embedding | embedding | embedding |

<sos>    je    suis

# 2. Attention

## Attention Mechanism

$$\hat{y}_t = \text{Softmax}\left(W_y \tilde{s}_t + b_y\right)$$

$$\tanh\left(\begin{array}{|c|}\hline W_c \\\hline\end{array} \times \begin{array}{|c|}\hline v_t \\\hline\end{array}\right) = \begin{array}{|c|}\hline \tilde{s}_t \\\hline\end{array}$$

**Attention Value $a_t$**

**concatenate**

$s_t$

$v_t$

softmax

| LSTM | → | LSTM | → | LSTM |

| embedding | embedding | embedding |

↑       ↑       ↑

&lt;sos&gt;      je      suis

# 2. Attention

## Attention Mechanism

| 이름 | 식 | 출처 |
|---|---|---|
| content-based attention | $f(s,\ h) = \dfrac{s^T h}{\lvert\lvert s \rvert\rvert \cdot \lvert\lvert h \rvert\rvert}$ | Graves, 2014 |
| additive attention (Bahdanau attention) | $f(s,\ h) = V^T \tanh(W_1 s + W_2 h)$[4] | Bahdanau, 2015 |
| dot-product attention (Loung attention) | $f(s,\ h) = s^T h$ | Luong, 2015 |
| scaled dot-product attention | $f(s,\ h) = \dfrac{s^T h}{\sqrt{n}}$ [5] | Vaswani, 2017 |

# 2. Attention

Attention Mechanism

But still,,,

1) Gradient Vanishing / Exploding

2) Too slow (non-parallel operation)

=> Transformer

# 3. Training Techniques

## Dropout



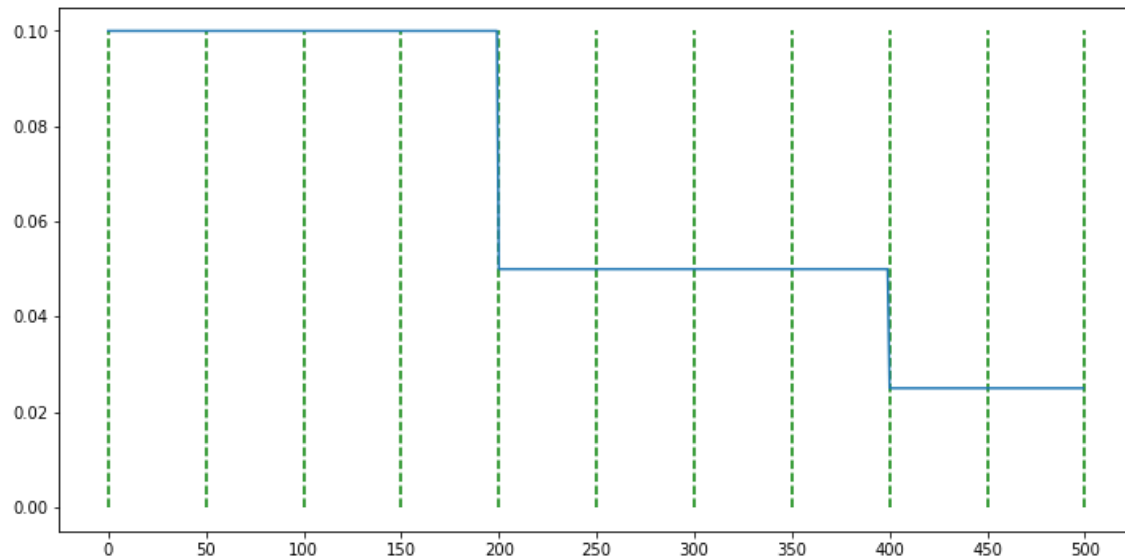(a) Standard Neural Net

(b) After applying dropout.

## LR Scheduler: Lambda LR Scheduler

```
scheduler = LambdaLR(optimizer, lr_lambda = lambda epoch: 0.95 ** epoch)
```
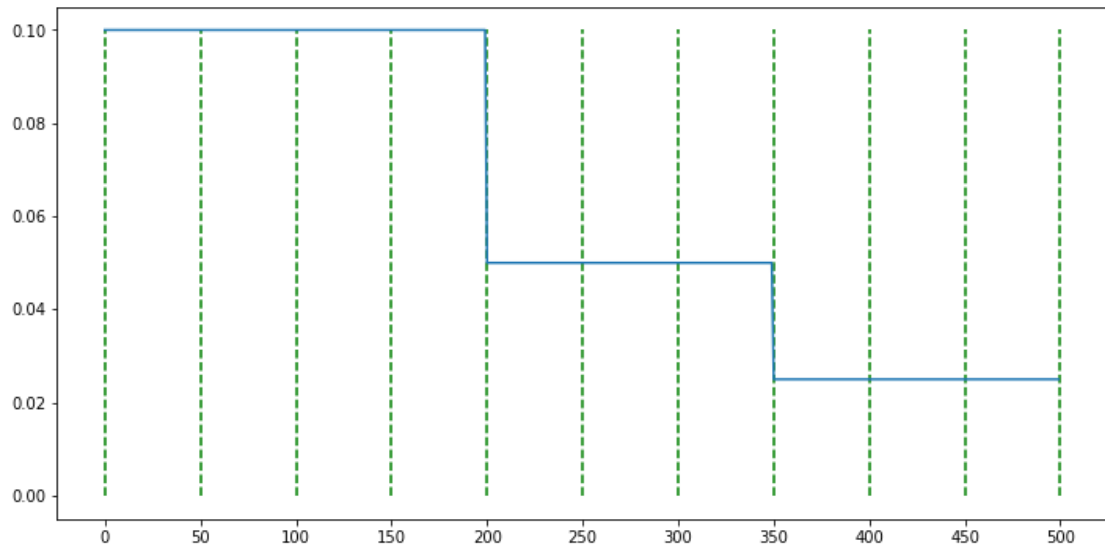
## Step LR Scheduler

```
scheduler = StepLR(optimizer, step_size=200, gamma=0.5)
```
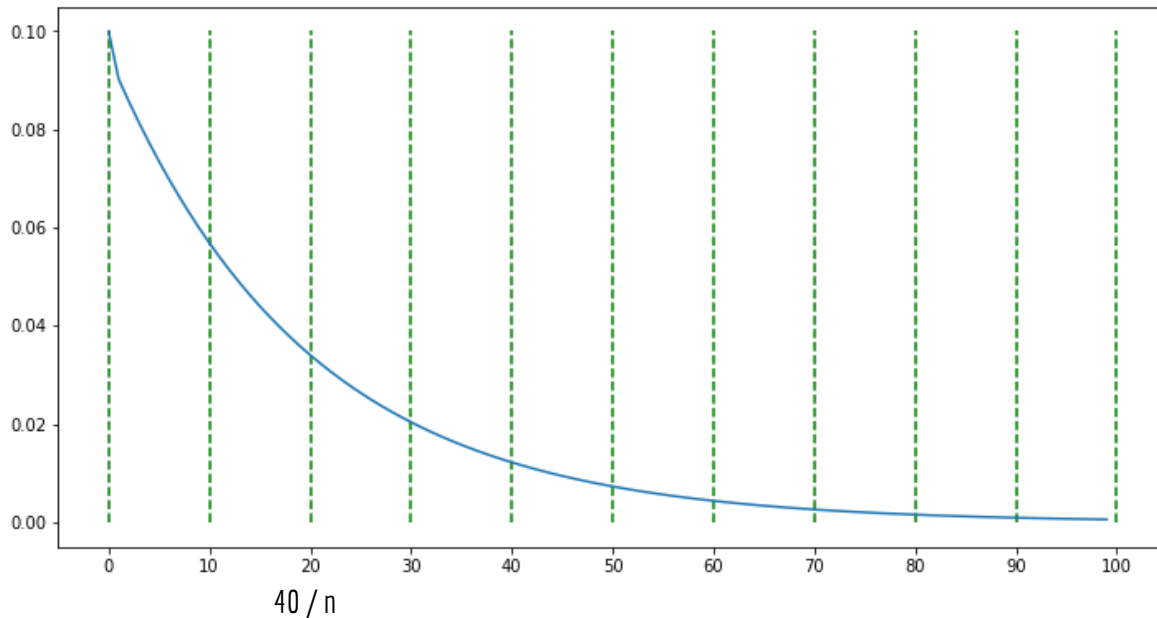
## Multi-Step LR Scheduler

```
scheduler = MultiStepLR(optimizer, milestones=[200, 350], gamma=0.5)
```

## Exponential LR Scheduler
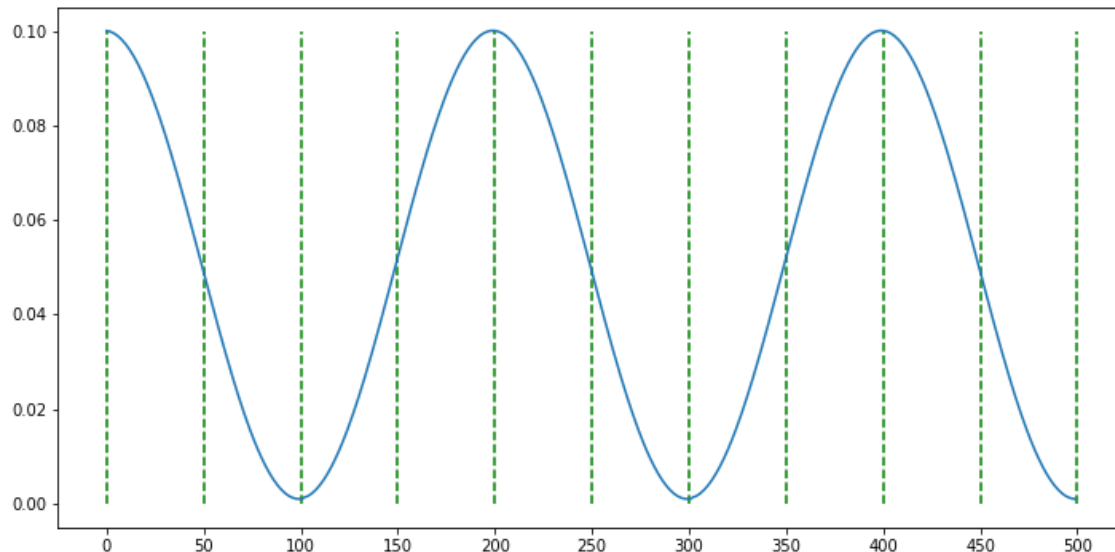
```
scheduler = ExponentialLR(optimizer, gamma=0.95)
```



40 / n

# 3. Training Techniques

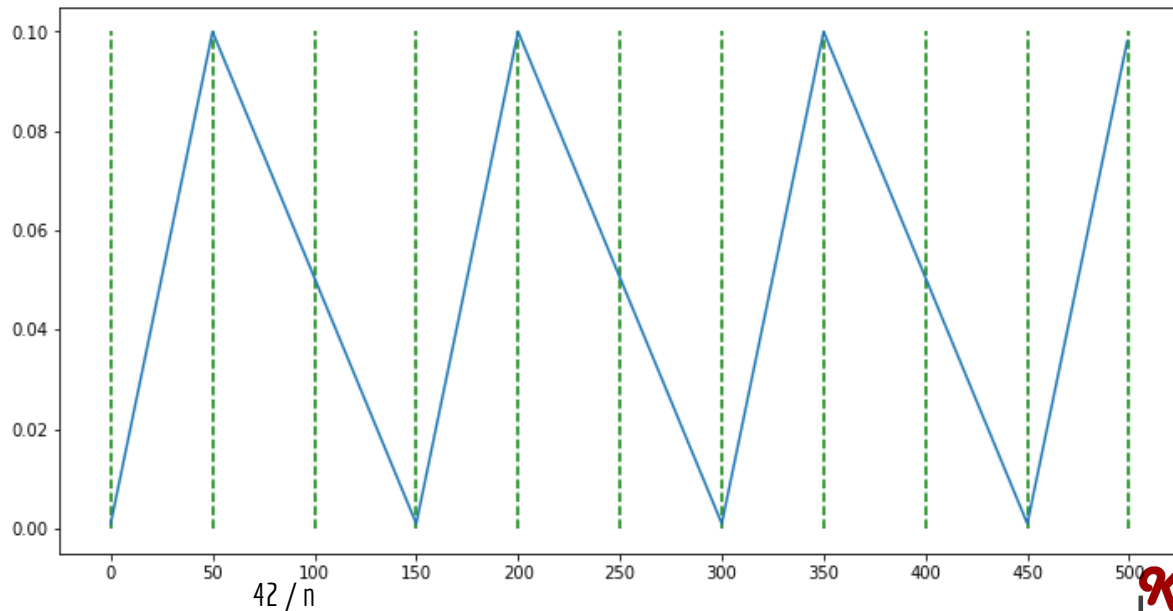## Cosine Annealing LR Scheduler

```
scheduler = CosineAnnealingLR(optimizer, T_max=100, eta_min=0.001)
```
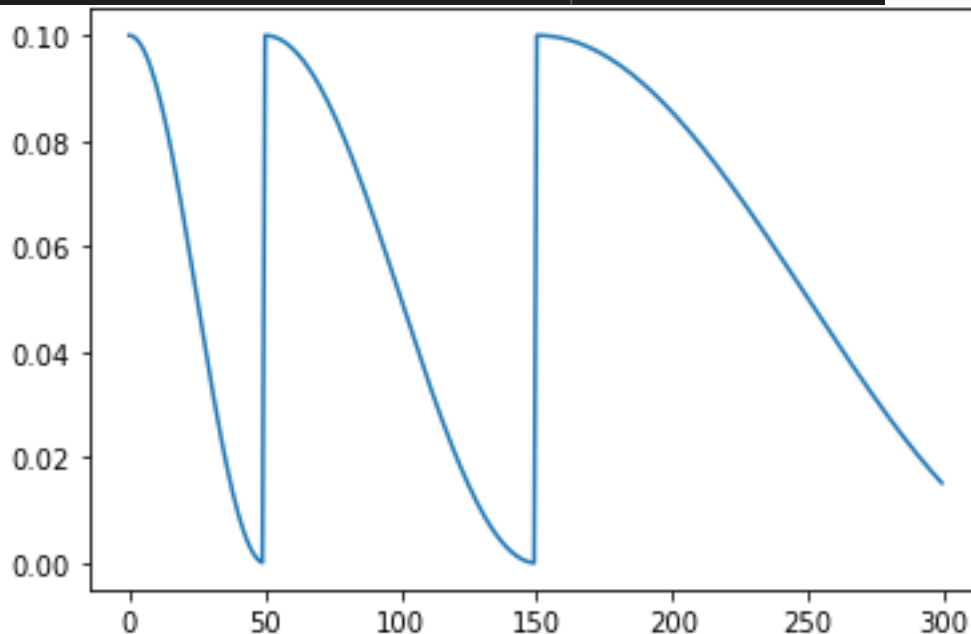
## Cycle LR Scheduler

```
scheduler = CyclicLR(optimizer, base_lr=0.001, max_lr=0.1, step_size_up=50, step_size_down=100, mode='triangular')
```
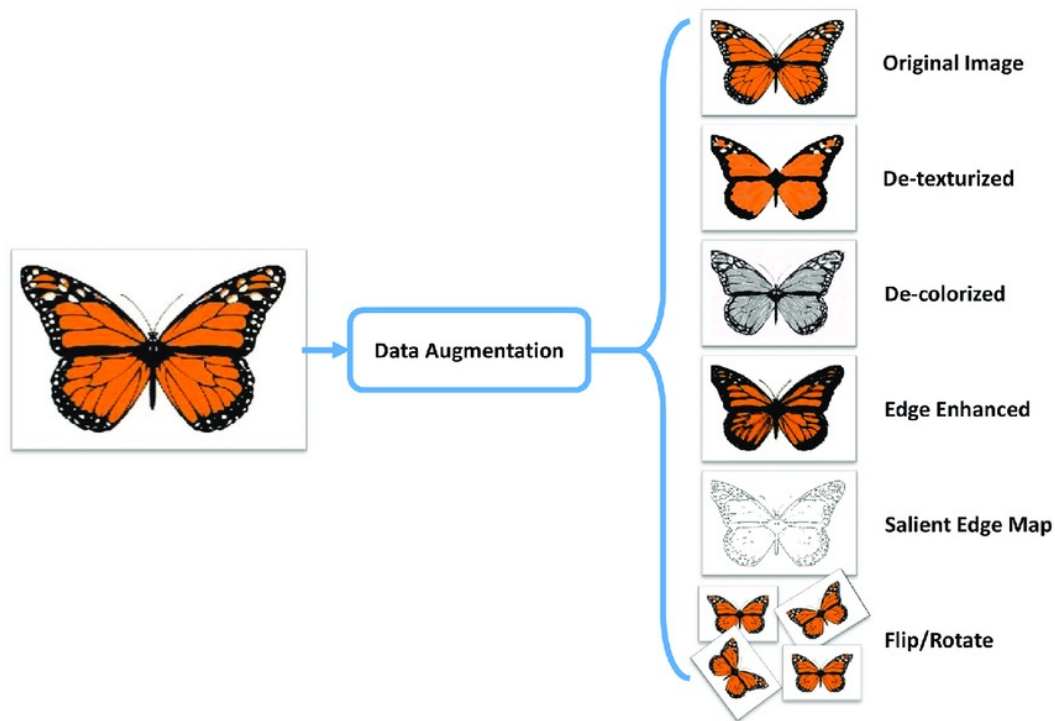


42 / n

# 3. Training Techniques

## Cosine Annealing LR Scheduler with Warm Restarts

```
scheduler = optim.lr_scheduler.CosineAnnealingWarmRestarts(optimizer, T_0=50, T_mult=2, eta_min=0)
```
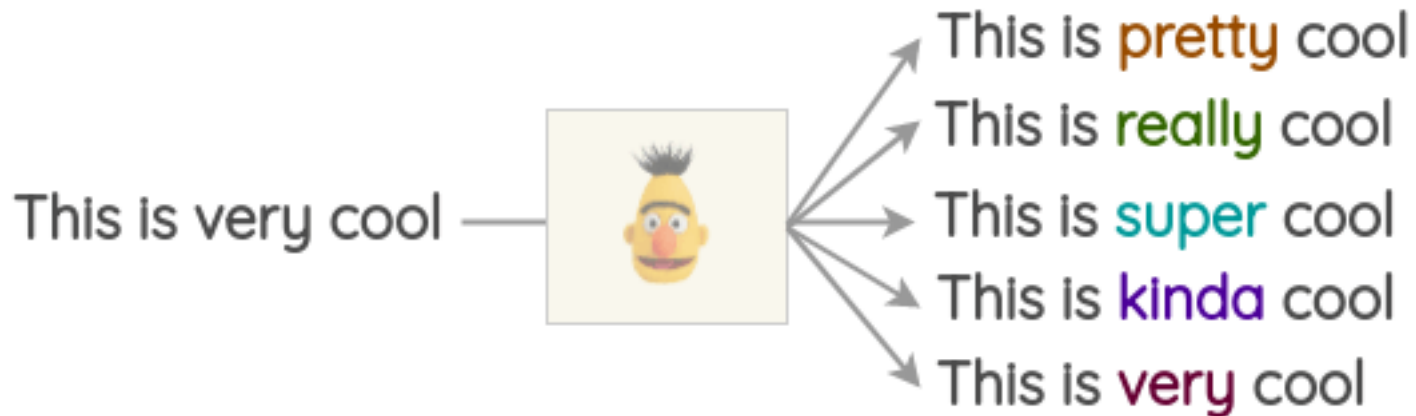
# 3. Training Techniques

## Data Augmentation

# 3. Training Techniques

## Data Augmentation



This is very cool → This is **pretty** cool
This is **really** cool
This is **super** cool
This is **kinda** cool
This is **very** cool

Does it seem efficient?

# Q&A

# Have a nice week🥰