**ICT이노베이션스퀘어 AI복합교육 고급 언어과정**

# 자연어처리를 위한
# Machine Translation
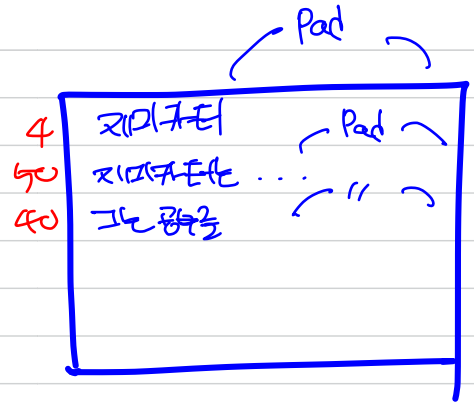
**현청천**

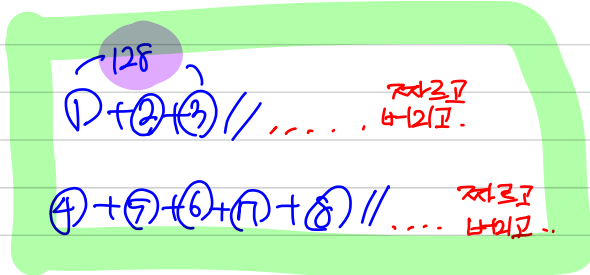2021.04.19

ex )

지미커터. 지미카터는 미국//대통령이다. 그는
<u>시민</u>(○)                      <u>시장</u> (x)

4+ 50+40+ 150

가는 도중에 자르기 (느낌)

| 4 | 지미커터        Pad |
| 50 | 지미카터는 ...  Pad |
| 40 | 그는 활동중 "   |

① ②
③ ④
⋮
⑩⑩ (100)

➡ [ ~128
①+②+③//... 제자고 나의고.

④+⑤+⑥+⑦+⑧//... 제자고 나의고. ]

○ 시작점은 같지만
○ 뒷쪽은 자르자.
○ but) 최대한 많은 문장들은 합쳐서
   긴 문장을 학습시킬수있는 장점이 있다.

이순신
이순신은...

이때 자르는 단위 ⇒ 정고 // = chunk
                    ⇒ 128 //

⇒ 문장대로 끊는게 더 안 좋다라는... 논문에서 ⇒ Robert

- Segment-Pair + NSP < 512  ⇒ 우리가쓰는거

∵ 단위문장보다
chunk로 만들어서 학습하는게
더 좋다.

— Sentence- Pair + NSP
  ○ 한줄 씩. 한 문장직만

- Full Sentence

지미카터
수학          ) 아예 큰 꼬꼬망 장황내자.

# What is Machine Translation

인간이 사용하는 <u>자연 언어</u>를 **컴퓨터**를 사용하여 <u>다른 언어</u>로 **번역**하는 것
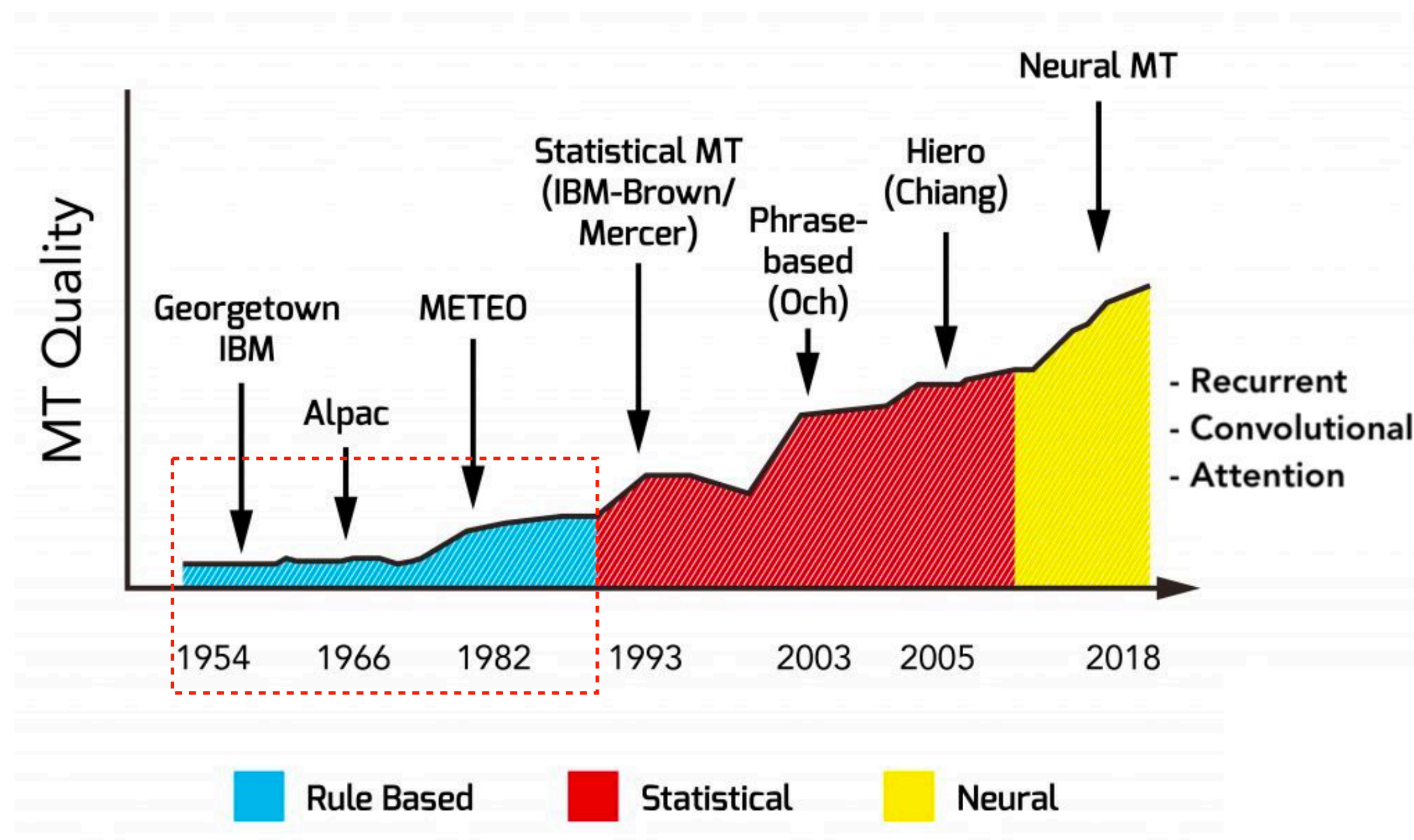
| Source (x) | Target (y) |
|---|---|

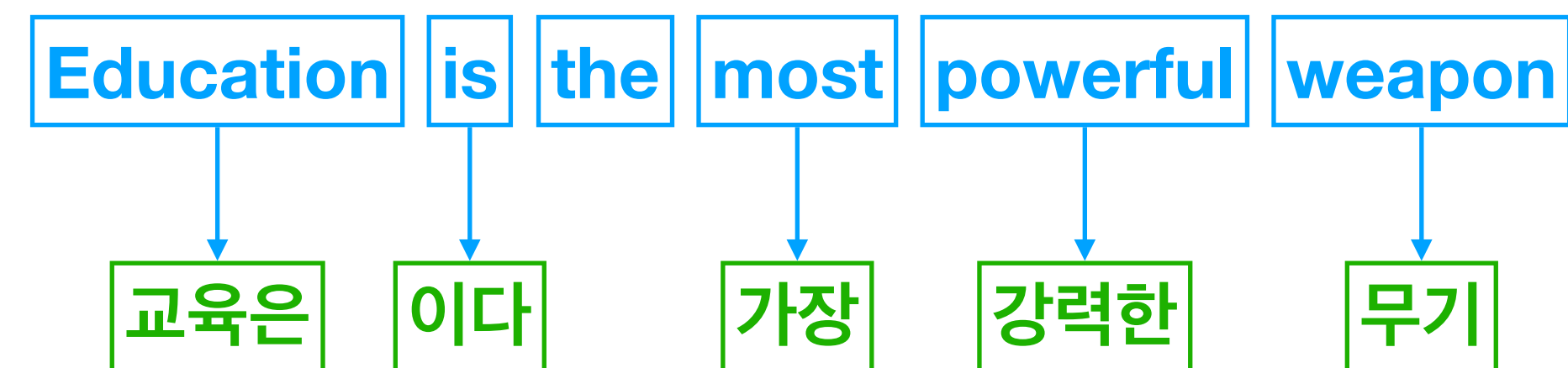**Education is the most powerful weapon we can use to change the world.**

**교육은 세상을 바꿀 수 있는 가장 강력한 무기이다.**

# What is Machine Translation (history)



**Rule-based Machine Translation**
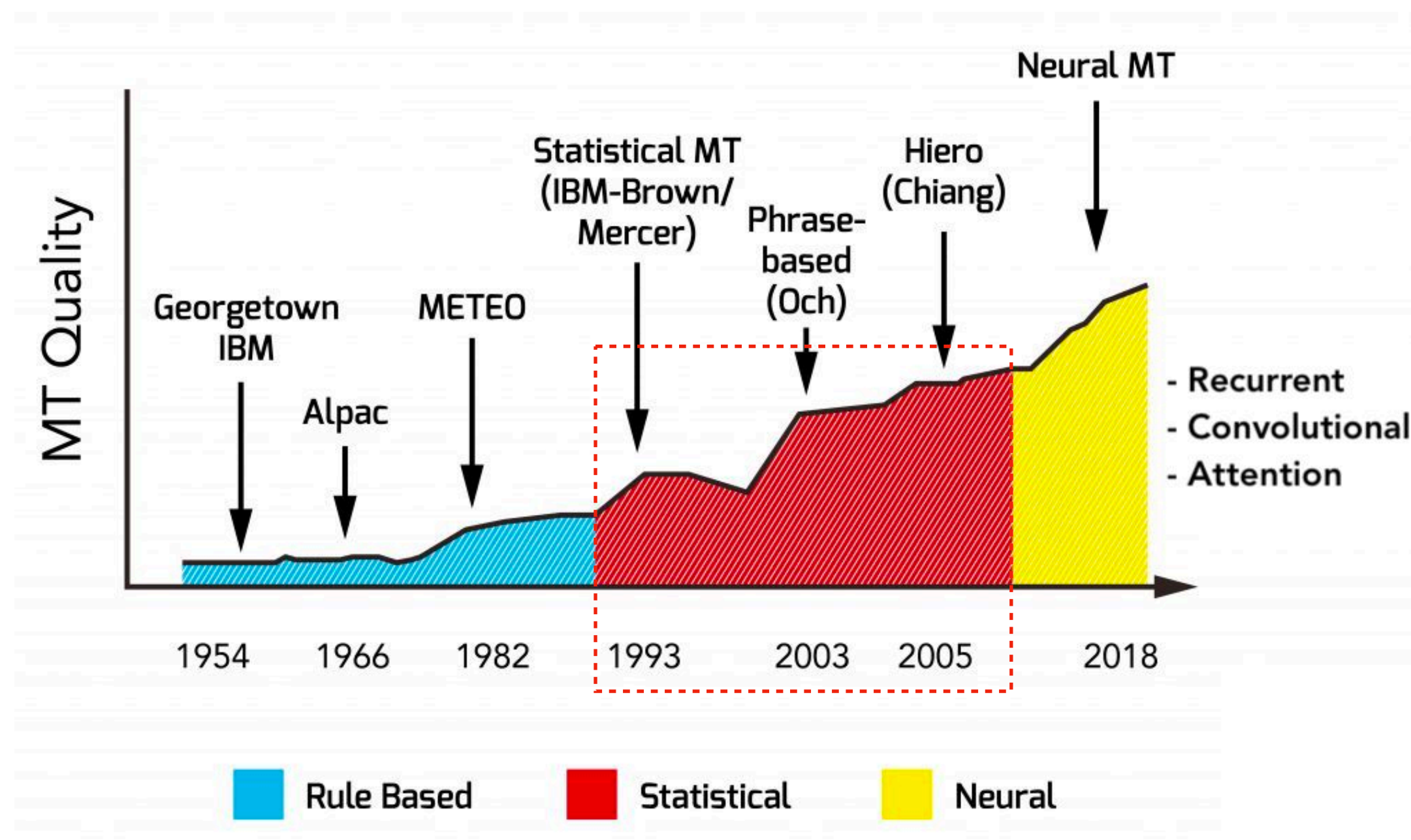
- Bilingual dictionary
- Linguistic rules for each language

| Education | is | the | most | powerful | weapon |
|-----------|-----|-----|------|----------|--------|
| 교육은 | 이다 | | 가장 | 강력한 | 무기 |

**I saw a man on a hill with a telescope?**

3

# What is Machine Translation (history)



## Statistical Machine Translation

- Language pair로부터 패턴 학습
- 데이터가 많을 수록 좋은 결과
- $argmax_y P(y|x)$

  ↖ input.

$x = $ 'Education is the most powerful weapon'

$$'교육은' = argmax_y P(y|x)$$

$$'가장' = argmax_y P(y|x, '교육은')$$

$$'강역한' = argmax_y P(y|x, '교육은', '가장')$$

- - - - - - - - -

이미지: https://iconictranslation.com/what-we-do/neural-machine-translation/
참조: https://www.freecodecamp.org/news/a-history-of-machine-translation-from-the-cold-war-to-deep-learning-f1d335ce8b5/

# What is Machine Translation (history)
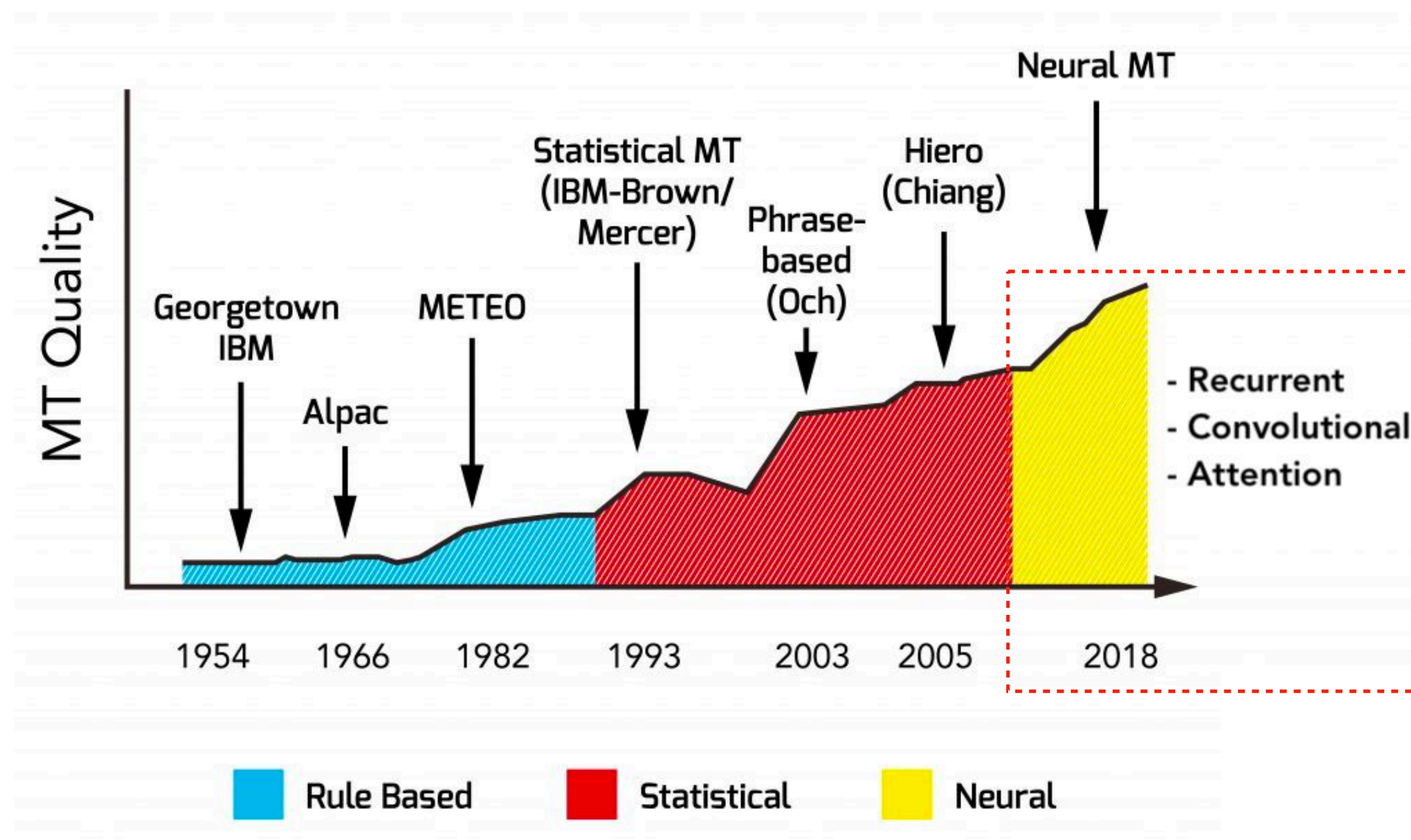


## Neural Machine Translation

- 2041년 sequence to sequence 모델 등장
- 데이터에서 Neural Network 학습
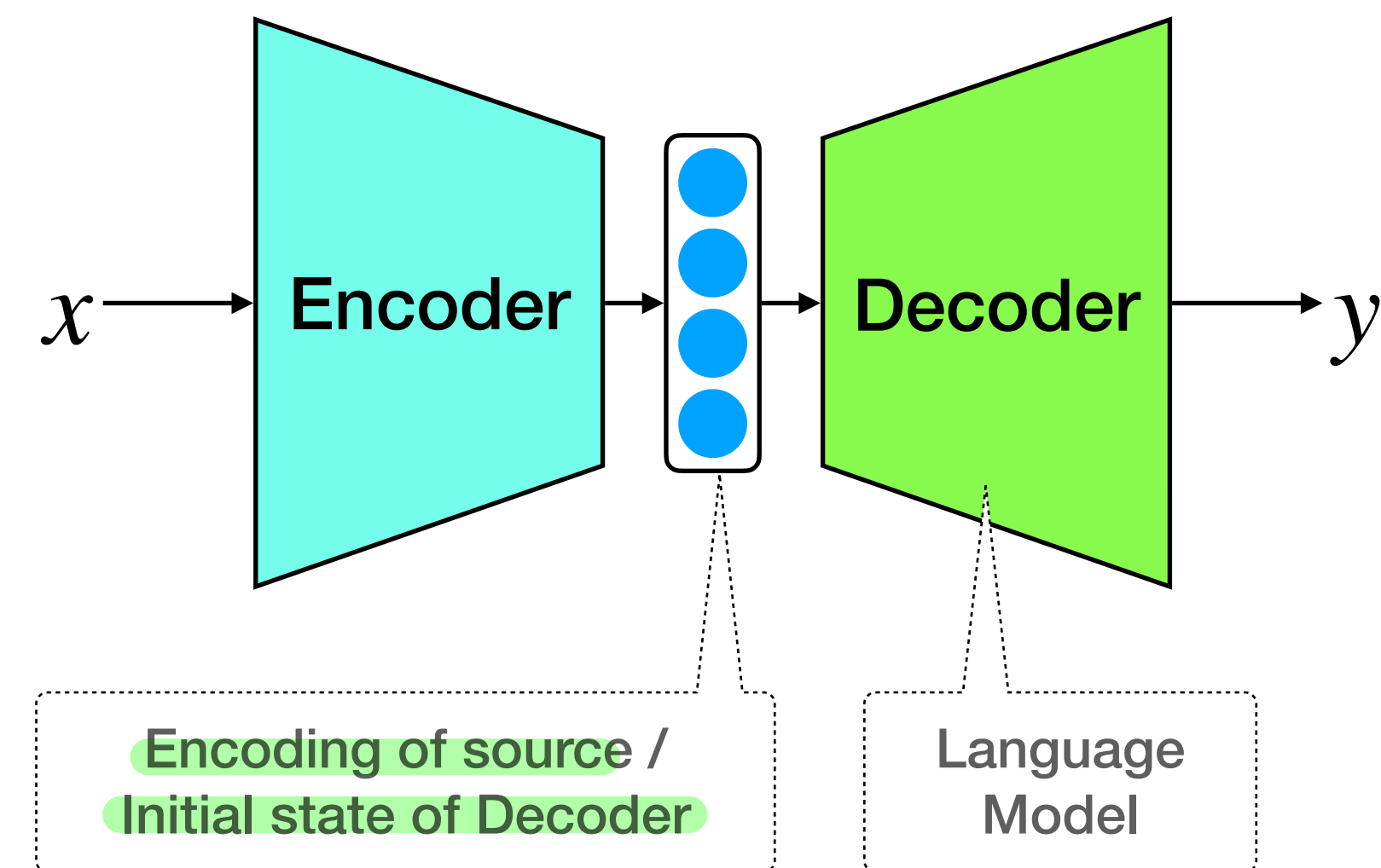- $P(y|x;\theta)$
  - 어순 오류 감소
  - 어휘 오류 감소
  - 문법 오류 감소

이미지: https://iconictranslation.com/what-we-do/neural-machine-translation/
참조: https://www.freecodecamp.org/news/a-history-of-machine-translation-from-the-cold-war-to-deep-learning-f1d335ce8b5/

# What is Machine Translation (history)



**Neural Machine Translation**

*sentence embedding*



$x \rightarrow$ Encoder $\rightarrow$ Decoder $\rightarrow y$

Encoding of source /
Initial state of Decoder

Language
Model

○ 디코더의 초기값
○ input을 Encoding 해서 Vector로 만들기!

이미지: https://iconictranslation.com/what-we-do/neural-machine-translation/
참조: https://www.freecodecamp.org/news/a-history-of-machine-translation-from-the-cold-war-to-deep-learning-f1d335ce8b5/

# Machine Translation DataSet

*ex) Bible,,은 분통 다 정체에 있나까..*

- WMT Dataset

*다 영어 기반되니..*

| File | CS-EN | DE-EN | IU-EN | JA-EN | KM-EN | PL-EN | PS-EN | RU-EN | TA-EN | ZH-EN | FR-DE |
|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Europarl v10 | ✓ | ✓ | | | | ✓ | | | | | ✓ |
| ParaCrawl v5.1 | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | | | ✓ |
| Common Crawl corpus | ✓ | ✓ | | | | | | ✓ | | | ✓ |
| News Commentary v15 | ✓ | ✓ | | ✓ | | | | ✓ | | ✓ | ✓ |
| CzEng 2.0 | ✓ | | | | | | | | | | |
| Yandex Corpus | | | | | | | | ✓ | | | |
| Wiki Titles v2 | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| UN Parallel Corpus V1.0 | | | | | | | | ✓ | | ✓ | |
| Tilde Rapid corpus | ✓ | ✓ | | | | ✓ | | | | | |
| CCMT Corpus | | | | | | | | | | ✓ | |
| WikiMatrix | ✓ | ✓ | | ✓ | | ✓ | | ✓ | ✓ | ✓ | ✓ |
| Back-translated news | ✓ | | | | | | | ✓ | | ✓ | |
| Japanese-English Subtitle Corpus | | | | ✓ | | | | | | | |
| The Kyoto Free Translation Task Corpus | | | | ✓ | | | | | | | |
| TED Talks | | | | ✓ | | | | | | | |
| Nunavut Hansard Inuktitut-English Parallel Corpus 3.0 | | | ✓ | | | | | | | | |
| PMIndia v1 | | | | | | | | | ✓ | | |
| Tanzil v1 | | | | | | | | | ✓ | | |

- Workshop on Statistical Machine Translation
- Bilingual Datasets
- English Based Datasets
- http://www.statmt.org/wmt20/translation-task.html

# Machine Translation DataSet

- WMT Dataset
- AI-Hub 한국어-영어 병렬 말뭉치

**한국어-영어 번역(병렬) 말뭉치 AI 데이터 다운로드**

| 소개 | 다운로드 | 저작도구 |

| 다운로드 ✏️ | 문의하기 ✏️ |

**한국어-영어 번역 말뭉치** [전체 선택] [선택 해제] [다운로드]

| ⬇ 구어체(1) 다운로드 | ⬇ 구어체(2) 다운로드 | ⬇ 대화체 다운로드 |
| ⬇ 문어체-뉴스(1) 다운로드 | ⬇ 문어체-뉴스(2) 다운로드 | ⬇ 문어체-뉴스(3) 다운로드 |
| ⬇ 문어체-뉴스(4) 다운로드 | ⬇ 문어체-한국문화 다운로드 | ⬇ 문어체-조례 다운로드 |
| ⬇ 문어체-지차체웹사이트 다운로드 | | |

- AI-Hub 한국어-영어 데이터
- 회원 가입 및 별도의 서류제출 후 다운로드
- https://aihub.or.kr/aidata/87/download

# Machine Translation Model

ex) convolutional Seq 2 Seq

○ 왜 학습이 Encoder → Decoder 형태로 만들어질까?

○ 왜 하나의 Vector로 만들어져야 되는가?

나는 학생 입니다

3개의 Vector

How? RNN, LSTM 등 가장 feature가 잘 뽑힌거로.

1개의 Vector로 만들자!

**Encoder**

**Decoder**

Target

Source

○ 번역할 문장

Encoding of source /
Initial state of Decoder

○ Sentence embedding

Language
Model

○ 생성

# Machine Translation Model

교육은 가장 강력한 무기이다

**Encoder** → **Decoder**

Education is the most powerful weapon

**Encoder Decoder Architecture / Sequence to Sequence**

# Machine Translation Model (Training)

Decoder Label
(Target + [EOS])

| 교육은 | 가장 | 강력한 | 무기 | 입니다 | [EOS] |

L.M

| Education | is | the | most | powerful | weapon |

Encoder Input
(Source)

| [BOS] | 교육은 | 가장 | 강력한 | 무기 | 입니다 |

Decoder Input
([BOS] + Target)

어학습할 모델

# Machine Translation Model (Training)

일반적으로 LSTM을 많이 사용함!

(encoder —hidden— state)
(encoder —cell —state)

Encoding of source /
Initial state of Decoder

Decoder Label
(Target + [EOS])

교육은   가장   강력한   무기   입니다   [EOS]

e2) RNN
LSTM H₁   H₂   ··   ··   ··   Hm

+

H₀에

Embedding

Education   is   the   most   powerful   weapon

Encoder Input
(Source)

[BOS]   교육은   가장   강력한   무기   입니다

Decoder Input
([BOS] + Target)

# Machine Translation Model (Training)



o Encoder (LSTM) → ① e-h
② e-h-state
③ e-c-state

oore그ue

(concat) × 2Hz

(bs.5)
(bs.5) → concat ⇒ (bs,10)
(bs.5) = `cat`이라할때

GRU (Unit =10, initial-state = [cat])

Decoder (GRU) → shape로 맞춰야할경우

**Decoder Label (Target + [EOS])**

교육은   가장   강력한   무기   입니다   [EOS]

Dense. softmax

Encoding of source / Initial state of Decoder

$h_m = S_0$

$S_1$   $S_2$   $S_3$   ...   $S_m$

Education | is | the | most | powerful | weapon

**Encoder Input (Source)**

[BOS]   교육은   가장   강력한   무기   입니다

**Decoder Input ([BOS] + Target)**

| En RNN | | D ~ RNN |
|---|---|---|
| Embedding | | Embedding |

↑　　　　　　　　　↑
영어　　　　　　　한글

○ 영어, 한글 일단 나눠쓰는게 일반적

○ but) 한글. 한글일때 같는걸 사용!

그렇지 않으면 다른걸로 까ₗ되어서 학습됩수있지.

# Machine Translation Model (Training)

Decoder Label
(Target + [EOS])

Encoding of source /
Initial state of Decoder

| 교육은 | 가장 | 강력한 | 무기 | 입니다 | [EOS] |

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} J^t(\theta)$$

**Mean of
negative log prob**

| Education | is | the | most | powerful | weapon |

| [BOS] | 교육은 | 가장 | 강력한 | 무기 | 입니다 |

Encoder Input
(Source)

Decoder Input
([BOS] + Target)

# Machine Translation Model (Training)

$$p(y_1, \ldots, y_n | x_1, \ldots, x_m) = \prod_{t=1}^{n} p(y_t | v, y_1, \ldots, y_{t-1})$$

# Machine Translation Model (Training)

$$p(y_1, \ldots, y_n | x_1, \ldots, x_m) = \prod_{t=1}^{n} p(y_t | v, y_1, \ldots, y_{t-1})$$

label

| 교육은 | 가장 | 강력한 | 무기 | 입니다 | [EOS] |



| [BOS] | 교육은 | 가장 | 강력한 | 무기 | 입니다 |

| Education | is | the | most | powerful | weapon |

input

input

$$\boxed{E} \quad - \quad \boxed{D}$$

① 한글⑨ Answer ③ Summary

④ 긍정 or 부정

[BOS]

① 영어 (번역)

② Question (챗봇)

③ 긴 문장 (요약)

④ 문장

model predict $\quad ( en-input \quad , \quad dec-input )$

$[q-id] \qquad [q-id] \qquad$ batch-size = 4 이니까

② 챗봇

go-backward = True

④ 한단계 token이
    다음의 Pad라
    다예(3 답(예1)
    나씨의
    의미

False.

원정  강씨  번째  Pad  Pad

[bos] [나도] [반가워] [Pad] [Pad]

① ←← 인코더에서는 가능한가능

② 디코더에서는 ←← 순서안된다

○ 원래라는 Pad가 뒤에 있을려야 go-backward가 줄고

○ 앞에서 Pad를 생성받으면 go-backward = False 가줄고

concat([forward-hidden, back-hidden])
concat([forward-cell, back-cell])

back-hidden
back-cell

forward-hidden
forward-cell

연정   앞에서   변기된  [Pad]  [Pad]          [BOS]

○ shape [bs, d-model*2], [x1], [x1], [x1], [x1]
    ○ concat 된거 아에
       줄거빈은

○ encoder는    임방향이 가능

○ decoder는    단방향 만. (forward-)
       -go-back하면 답보고 차는 꼴이 되짐나.

feature

여기서 반복.

E → [////] → D

    ○ 문장이 길수마다 E에서 feature를 뽑을거
    ○ but) exp에서는 짧은 문장. 간단한 것이나

D의 input이 새로 생성될거이따.
E를 다시, 매번 feature 뽑을순봄!

# Machine Translation Model (Inference)

Education | is | the | most | powerful | weapon

Encoder Input
(Source)

[BOS]

Decoder Input
([BOS] + Target)

# Machine Translation Model (Inference)

Encoding of source /
Initial state of Decoder

RNN
LSTM ELUGI



| Education | is | the | most | powerful | weapon |
|-----------|-----|-----|------|----------|--------|

Encoder Input
(Source)

[BOS]

Decoder Input
([BOS] + Target)

# Machine Translation Model (Inference)

greedy
random〈P〉

Encoding of source /
Initial state of Decoder

교육은

*Sample*

Education | is | the | most | powerful | weapon

Encoder Input
(Source)

[BOS]

Decoder Input
([BOS] + Target)

(bs) 1 [3,5,4]
(bs=1) ㅁ-seq

d-model
0
1
2
3
4
5
6

게네기능

[3.000] [5,0,0,0] [4.0.0.0]

[1, 3, d-model]

# Machine Translation Model (Inference)

Encoding of source /
Initial state of Decoder

*Sample*

가장

Education    is    the    most    powerful    weapon

[BOS]    교육은

Decoder Input
([BOS] + Target)

Append sampled token

# Machine Translation Model (Inference)

Encoding of source /
Initial state of Decoder

강력한

*Sample*

Education | is | the | most | powerful | weapon

Encoder Input
(Source)

[BOS] | 교육은 | 가장

Decoder Input
([BOS] + Target)

Append sampled token

# **Machine Translation Model (Inference)**

Encoding of source /
Initial state of Decoder

Sample

무기

Education | is | the | most | powerful | weapon

[BOS] | 교육은 | 가장 | 강력한

Decoder Input
([BOS] + Target)

Append sampled token

22

# Machine Translation Model (Inference)

Encoding of source /
Initial state of Decoder

EOS 가 최대길이 까지 생성
(64)

입니다

*Sample*

| | | | | | |
|---|---|---|---|---|---|
| Education | is | the | most | powerful | weapon |

[BOS]  교육은  가장  강력한  무기

Encoder Input
(Source)

이 문장에 대한 번역 ↓

Decoder Input
([BOS] + Target)

Append sampled token

23

# **Machine Translation Model (Inference)**

CNN, Lstm, maxpooing 등b

feature을 넘하되 의미를 두자!

End Inference

Encoding of source /
Initial state of Decoder

[EOS]

*Sample*



| Education | is | the | most | powerful | weapon |

| [BOS] | 교육은 | 가장 | 강력한 | 무기 | 입니다 |

Encoder Input
(Source)

Decoder Input
([BOS] + Target)

Append sampled token

24

# Machine Translation Model (Versatile)

- Summarization (long text → short text)
- Dialog (user utterance → agent utterance)
- Parsing (text → parsed sequence)
- Code generation (text → program code)
- etc.

# Machine Translation Model (Decoding)

① greedy search의 문제점..

교육은 가장 30,5% 30%

| 교육은 | 가장 | 강력한 | 무기 | 입니다 | [EOS] |

P 확률값   argmax   argmax   argmax   argmax   argmax   argmax

[BOS]   교육은   가장   강력한   무기   입니다

쿼리

argmax( ) 3

**Greedy Search**

Top 1 만

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$

즈 · $P_{L/m}$

[BOS]

$(k=5)$ 까지 함

$\varsigma$

ㅇ 문장을 고개요 뽑고
그중에 좋은 걸고 해1~

**Beam Search (k=2)**

# Machine Translation Model (Decoding)

$\sum_{i=1}^{\pm} P \Rightarrow P$ 줄이기

$i=1$

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$

더하기

$-0.7 = \log P_{LM}(\text{교육은} | [\text{BOS}])$

Top 2

교육은

[BOS]

운동은

$-0.9 = \log P_{LM}(\text{운동은} | [\text{BOS}])$

**Beam Search (k=2)**

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$

$-1$

$-1.7 = \log P_{LM}(\text{가장} | \text{[BOS], 교육은}) - 0.7$

*Top 2*

가장

여전 P
대비교

$-0.7$

교육은

많은

$-2.2$

$-2.9 = \log P_{LM}(\text{많은} | \text{[BOS], 교육은}) - 0.7$

[BOS]

$-1.6 = \log P_{LM}(\text{좋은} | \text{[BOS], 운동은}) - 0.9$

*Top 2*

좋은

$-0.7$

운동은

즐거운

$-0.9$

$-0.9$

$-1.8 = \log P_{LM}(\text{즐거운} | \text{[BOS], 운동은}) - 0.9$

## Beam Search (k=2)

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$



**Beam Search (k=2)**

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i \,|\, y_1, \ldots, y_{i-1})$$

$-2.9 = \log P_{LM}(\text{빠른} \,|\, [\text{BOS}], \text{교육은}, \text{가장}) - 1.7$

빠른

$w_1, w_2$  가장

가장

$-1.7$

$-0.7$

교육은

강력한

많은

$-2.9$

$-2.5 = \log P_{LM}(\text{강력한} \,|\, [\text{BOS}], \text{교육은}, \text{가장}) - 1.7$

$-2.8 = \log P_{LM}(\text{환경이} \,|\, [\text{BOS}], \text{운동은}, \text{좋은}) - 1.6$

[BOS]

환경이

$-1.6$

좋은

운동은

생각이

$-0.9$

즐거운

$-1.8$

$-3.8 = \log P_{LM}(\text{생각이} \,|\, [\text{BOS}], \text{운동은}, \text{좋은}) - 1.6$

## Beam Search (k=2)

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i \mid y_1, \ldots, y_{i-1})$$



**Beam Search (k=2)**

# Machine Translation Model (Decoding)

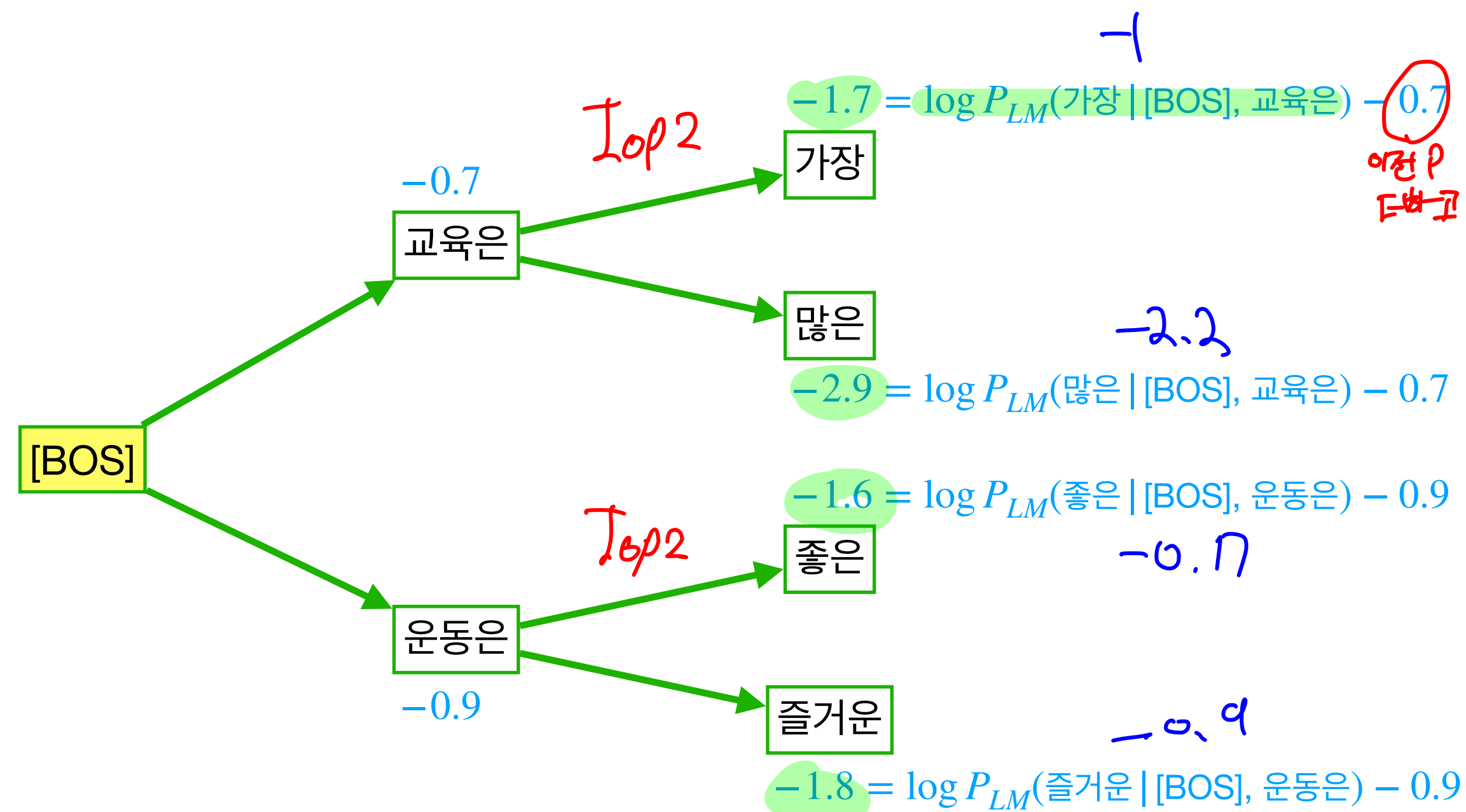$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i \mid y_1, \ldots, y_{i-1})$$



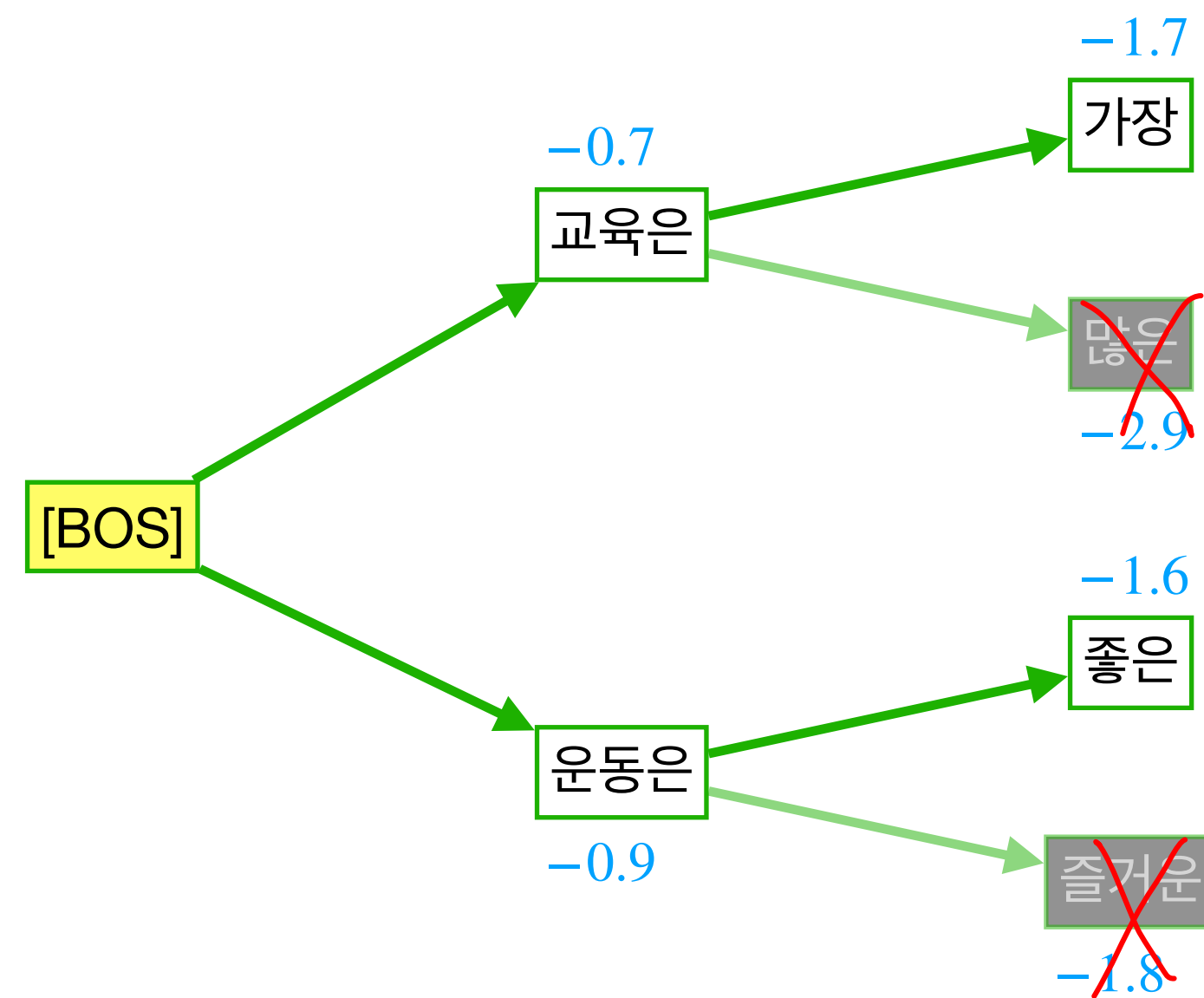**Beam Search (k=2)**

# Machine Translation Model (Decoding)

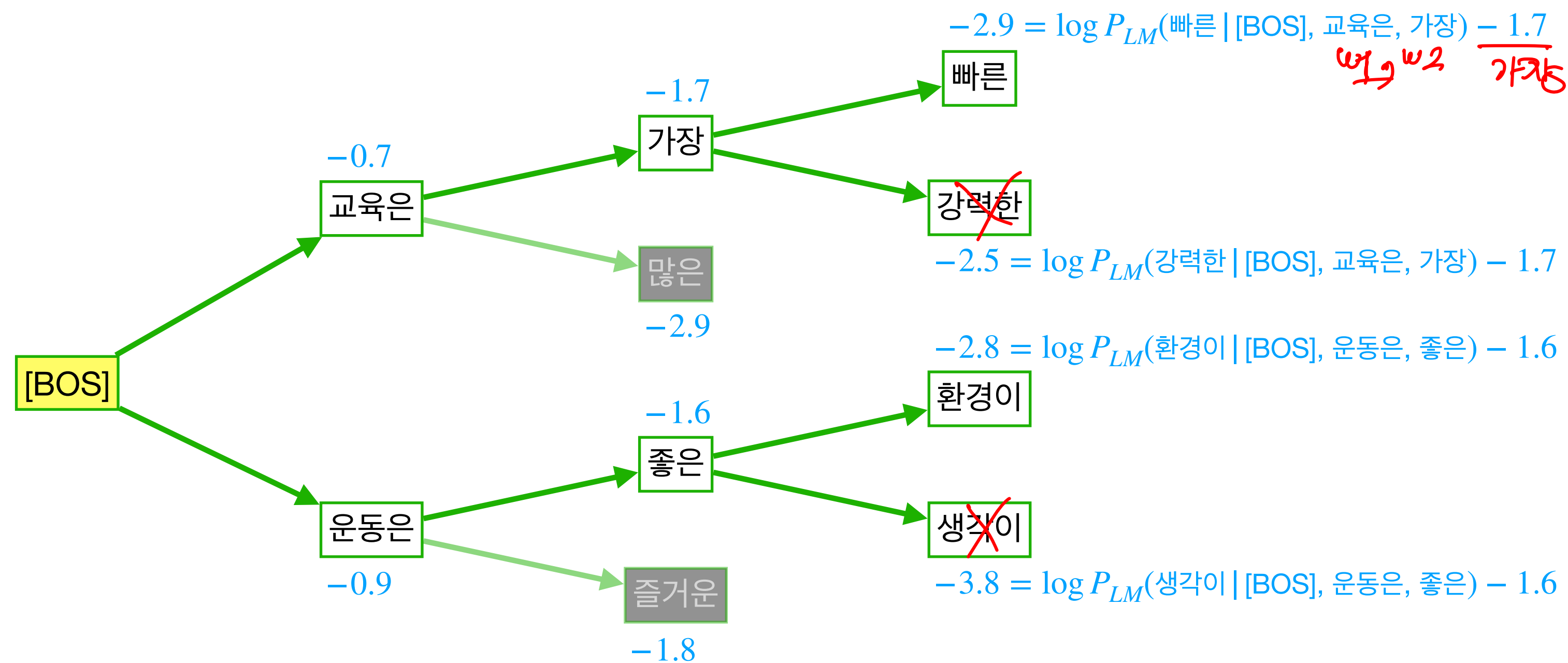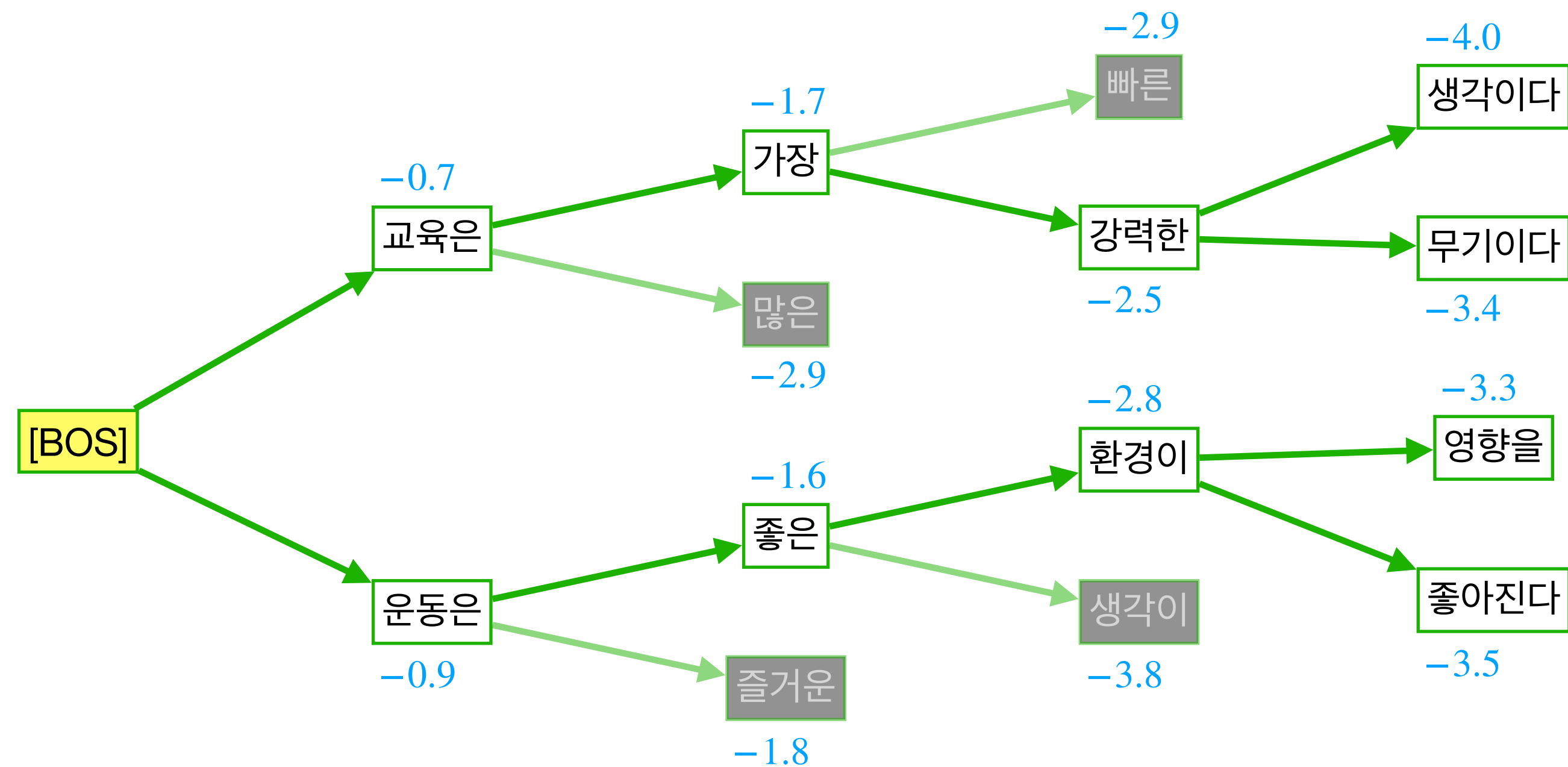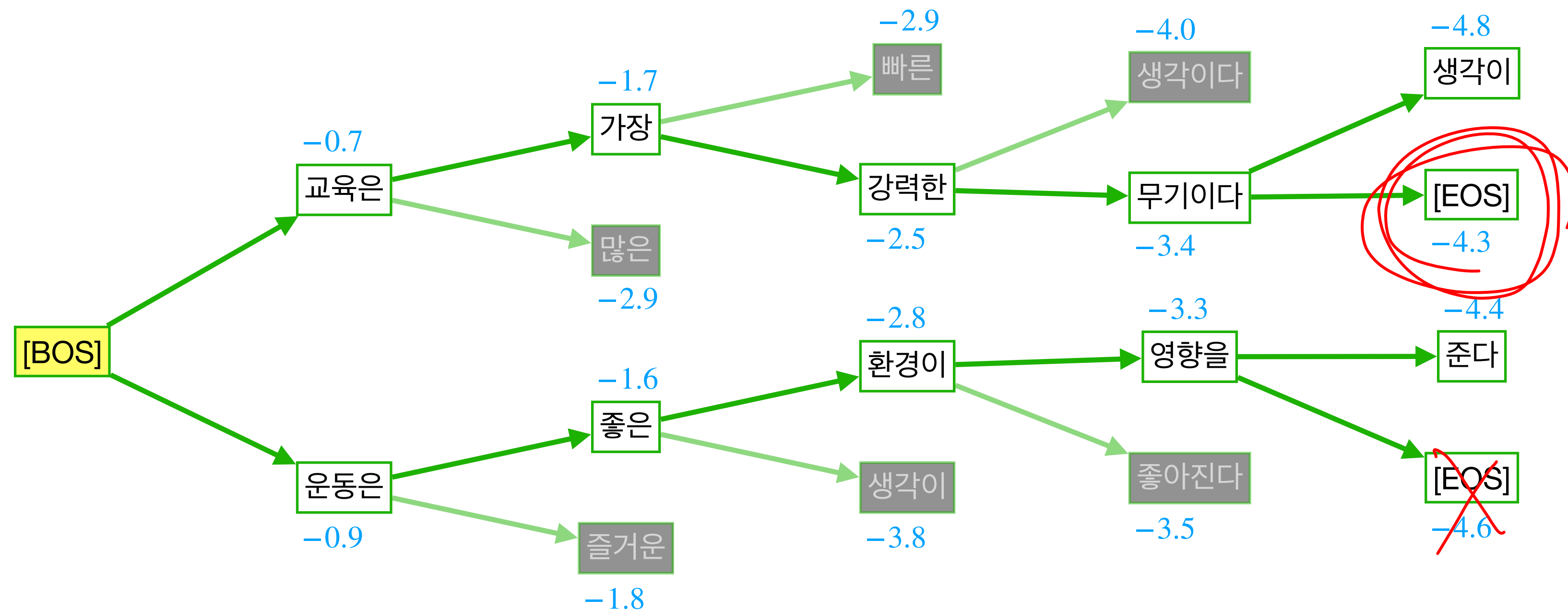$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$



**Beam Search (k=2)**

# Machine Translation Model (Decoding)

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$



**Beam Search (k=2)**

# Machine Translation Model (Decoding)

이 -1 아까 $-\infty$ . 까지나옴!

32001까지 다분류문제니까

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{i-1})$$

Highest Probability

−2.9
빠른

−4.0
생각이다

−4.8
생각이

−1.7
가장

−0.7
교육은

강력한

무기이다

[EOS]

−4.3 //

많은

−2.5

−3.4

−4.3

−5.6
많이

[BOS]

−2.9

−2.8
환경이

−3.3
영향을

−4.4
준다

−1.6
좋은

[EOS]

운동은

생각이

좋아진다

[EOS]

−5.1

−0.9

즐거운

−3.8

−3.5

−4.6

−1.8

−5.1

이 문장이 길수록 score가 낮아진다.

## Beam Search (k=2)

# Machine Translation Model (Decoding)

## Length Penalty

- Longer hypotheses have lower score  길어질수록 태하지네가 Score가 더 낮아지는거지 -

$$score(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{n-1})$$

=) 값이 작게게 뭐우건 줄겠냐

normalize 하중자능

- Normalize by length

$$score(y_1, \ldots, y_t) = \frac{1}{t} \sum_{i=1}^{t} \log P_{LM}(y_i | y_1, \ldots, y_{n-1})$$

단어의 길이만큼
평균 하우면  해결 하준가니!

## Beam Search (k=2)

# Machine Translation Model (Metric)

Education is the most powerful weapon we can use to change the world.

⇒ Q 어떻게 평가할건데요

## How do you evaluate???

교육은 세상을 바꿀 수 있는 가장 강력한 무기이다. ✓

세상을 바꿀 수 있는 가장 강력한 무기는 교육이다. ✓

가장 강력한 무기인 교육을 통해 세상을 바꿀 수 있다. ✓

뭐가 좋은지 알수 없잖아!

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

Candidate

reference

- Machine Translation된 결과와 사람이 Translation한 결과를 비교하여 품질을 평가

  True

  - n-gram precision ) n-gram 단위로 .
  - penalty for too-short system translations ) 짧은 문장에는 페널티 준다

$$\frac{TP}{TP + FP} \Rightarrow \frac{실제 \ 정답.}{정답하고 \ 예측한건중에}$$

# Machine Translation Model (Metric)

**BLEU(Bilingual Evaluate Understudy) Score**

**N-gram precision**

Candidate 1: It is a guide to action which ensures that the military always obeys the commands of the party.

Candidate 2: It is to insure the troops forever hearing the activity guidebook that party direct.

Reference 1: It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2: It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3: It is the practical guide for the army always to heed the directions of the party.

$$\text{Candidate 1 Unigram Precision} = \frac{17}{18}$$

$$\text{Candidate 2 Unigram Precision} = \frac{8}{14}$$

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

### N-gram precision

Candidate: the the the the the the the.

Reference 1: The cat is on the mat. =2개

Reference 2: There is a cat on the mat. =1개

> max =) 2

$$Count_{clip} = \min(Count, Max\_Ref\_Count)$$

2개

$$Candidate\ Unigram\ Precision = \frac{7}{7}$$

$$Candidate\ Modified\ Unigram\ Precision = \frac{2}{7}$$

2 = (1, 2)

# Machine Translation Model (Metric)

**BLEU(Bilingual Evaluate Understudy) Score**

**N-gram precision**

1
2
3
4
.
5

Candidate의 n-gram의 $Count_{clip}$개수

$$p_n = \frac{\sum_{n\text{-}gram \in C} Count_{clip}(n\text{-}gram)}{\sum_{n\text{-}gram' \in C} Count(n\text{-}gram')}$$

Candidate의 n-gram 개수

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

**Penalty for too-short system translations**

Candidate: of the  *(handwritten: 짧게 예측.)*

*(handwritten: 0 사실상  0 점짜리  번역인지  => 근데  blue가 만점,)*

Reference 1: It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2: It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3: It is the practical guide for the army always to heed the directions of the party.

$$\text{Candidate Unigram Precision} = \frac{2}{2}$$

*(handwritten: 당연히  Precision이  높아지지 .   Panelty를주자!)*

**짧은 문장일수록 n-gram precision이 높아지는 경향이 있음**

참고: https://www.aclweb.org/anthology/P02-1040.pdf

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

**Penalty for too-short system translations**

Candidate: of the

Reference 1: It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2: It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3: It is the practical guide for the army always to heed the directions of the party. $len \Rightarrow 16$

candidate > reference 들중 가장작은 len

brevity penalty

$$BP = \begin{cases} 1, & \text{if } c > r^{16} \\ exp(1 - \frac{r}{c}), & \text{if } c \leq r \end{cases}$$

= 음수여까.

ex) $exp(1 - \frac{16}{4})$

$= exp(-3)$

$< 1$

$BP = exp(1 - \frac{r}{c})$ 는 $exp < 1$ 이지 $\Rightarrow$ (페널티)

**짧은 문장일수록 n-gram precision이 높아지는 경향이 있음**

참고: https://www.aclweb.org/anthology/P02-1040.pdf

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

| gram | w | w |
|------|------|-------|
| 1 | 0.25 | 0.125 |
| 2 | 0.25 | 0.125 |
| 3 | 0.25 | 0.25 |
| 4 | 0.25 | 0.5 |

gram

0 또는 n=4, $w_n = \frac{1}{N}$

$= \frac{1}{4}$
$= 0.25$

$$p_n = \frac{\sum_{n\text{-}gram \in C} Count_{clip}(n\text{-}gram)}{\sum_{n\text{-}gram' \in C} Count(n\text{-}gram')}$$

$$BP = \begin{cases} 1, & \text{if } c > r \\ exp(1 - \frac{r}{c}), & \text{if } c \leq r \end{cases}$$

Sum

2-ip

$$BLEU = BP \cdot exp\left(\sum_{n=1}^{N} w_n \log p_n\right)$$

= blue Penalty

log ksn

**Baseline:** $N = 4, \quad w_n = \frac{1}{N}$

$$\log BLEU = min(1 - \frac{r}{c}, 0) + \sum_{n=1}^{N} w_n \log p_n$$

$(0.25 \quad 0.25 \quad 0.25 \quad 0.25)$

uni     bi     thr     half

# Machine Translation Model (Metric)

## BLEU(Bilingual Evaluate Understudy) Score

Candidate: 나는 어제 집에 가서 잠을 잤다

Reference: 나는 어제 집에 가서 잠을 설쳤다

*(handwritten) 번역 잘안했지?*

*(handwritten, red) Metric도 와 자꾸 count 기반을책계용*

*(handwritten, red) Metric은 연구하는 분야가 많이 있음.*

- BLEU는 유용한 지표지만 완벽하지 않음
- BLEU가 높으면 번역의 품질이 좋을 가능성이 높음
- 통계적인 지표

*(handwritten, blue) 그럼에도*

*(handwritten, red) 평가하는 방법이 마땅히 않다까.*

# 감사합니다.