



# 언어 모델 학습과 사용 쉽게 하세요!

## 효율적이고 확장성있는 사내 라이브러리 개발기

김정환 NAVER CLOVA  
안휘진 NAVER CLOVA

# CONTENTS

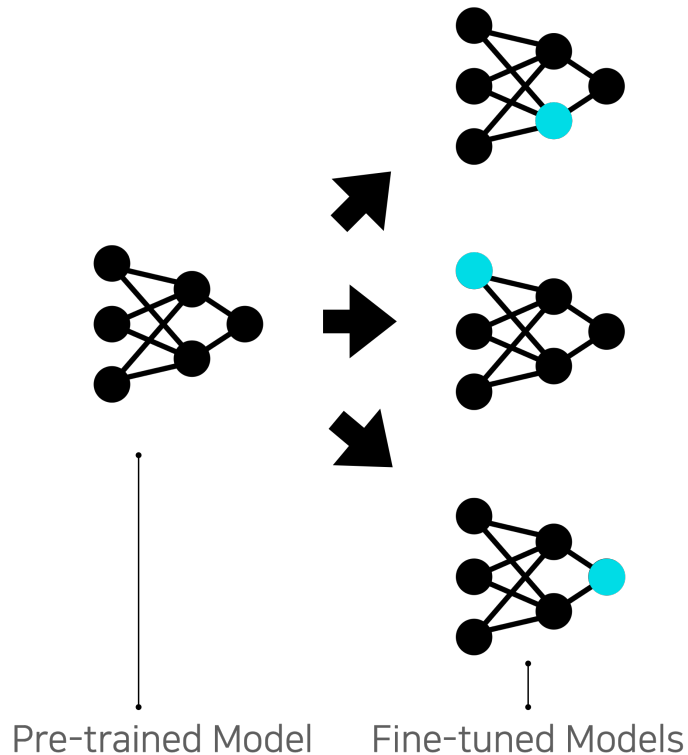
1. Pre-trained Language Model (PLM) 프로젝트 개발 배경
2. PPL: Pre-train Pipeline
3. LaRva: Language Representations by Clova
4. Future Work

# 1. 프로젝트 개발 배경

# Transfer Learning

## Fine-tuning

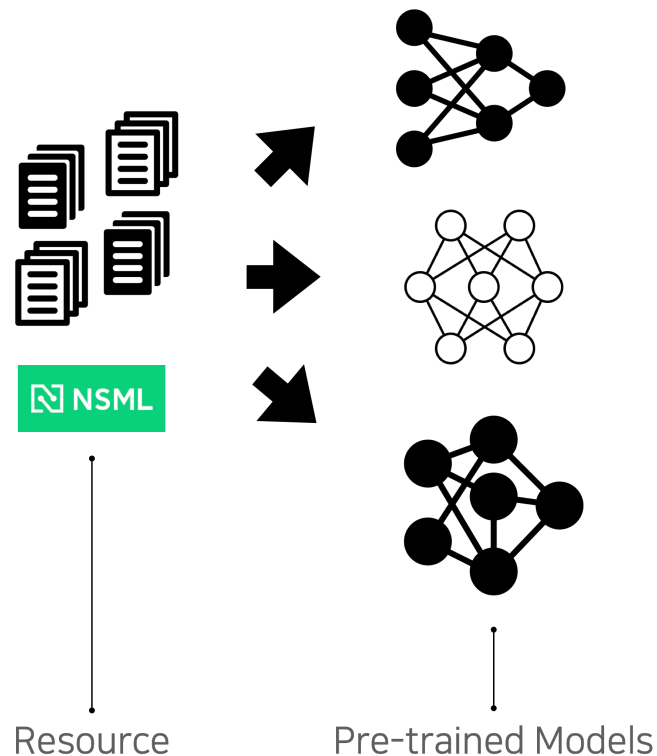
- Pre-train 단계에서 배운 지식의 전이
- 높은 성능
- 다양한 수요



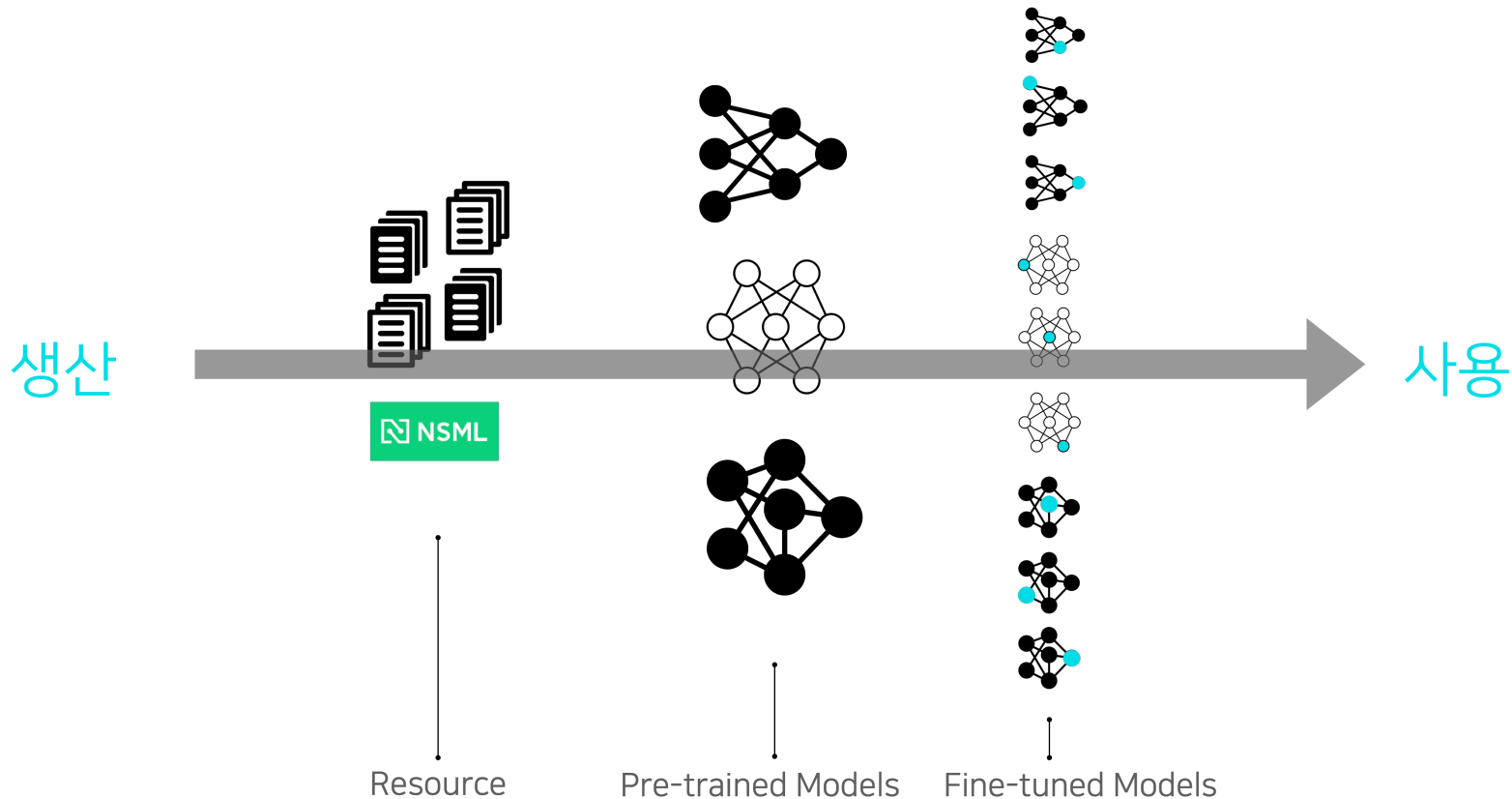
# Transfer Learning

## 다양한 Pre-trained Language Model (PLM)

- 사용된 코퍼스
- 모델 사이즈
- 토큰나이저
- ...



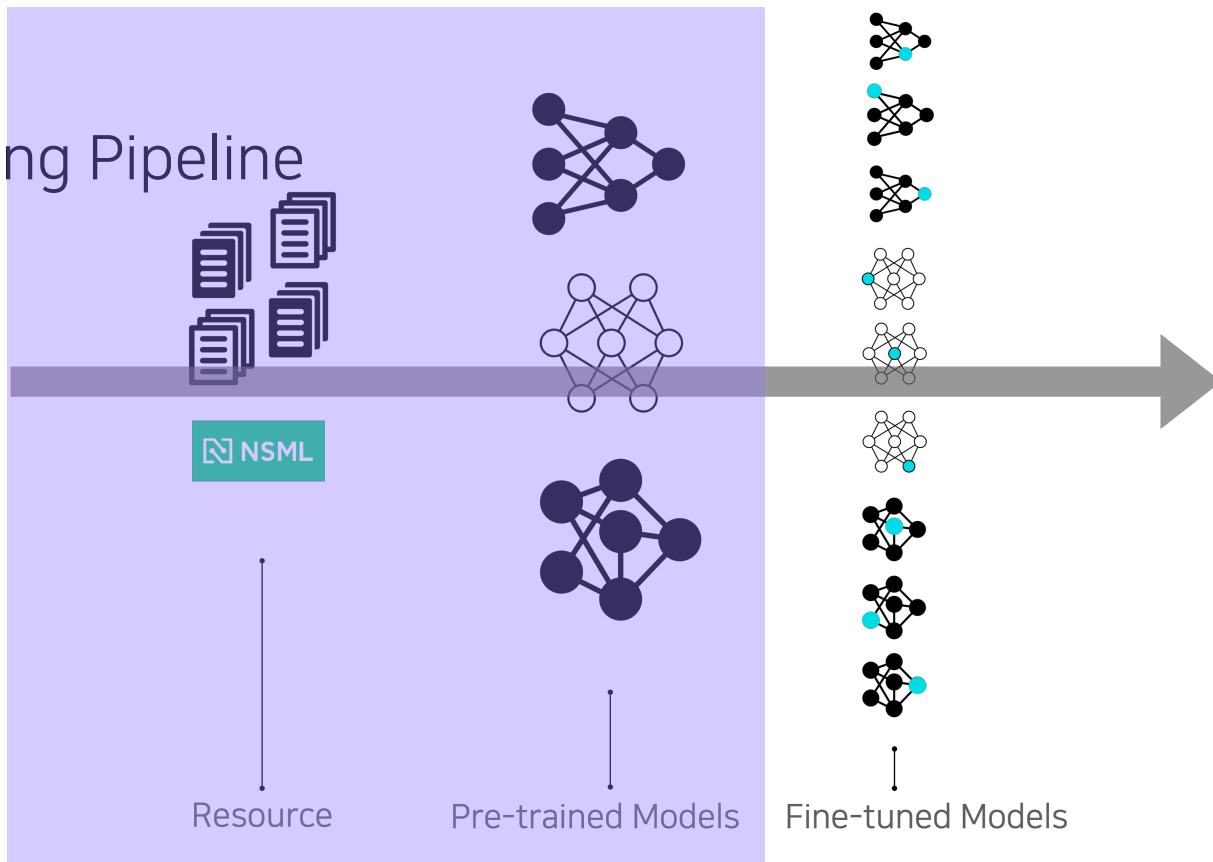
# PLM 생산과 사용 라이브러리



# PLM 생산과 사용 라이브러리

PPL:

Pre-training Pipeline

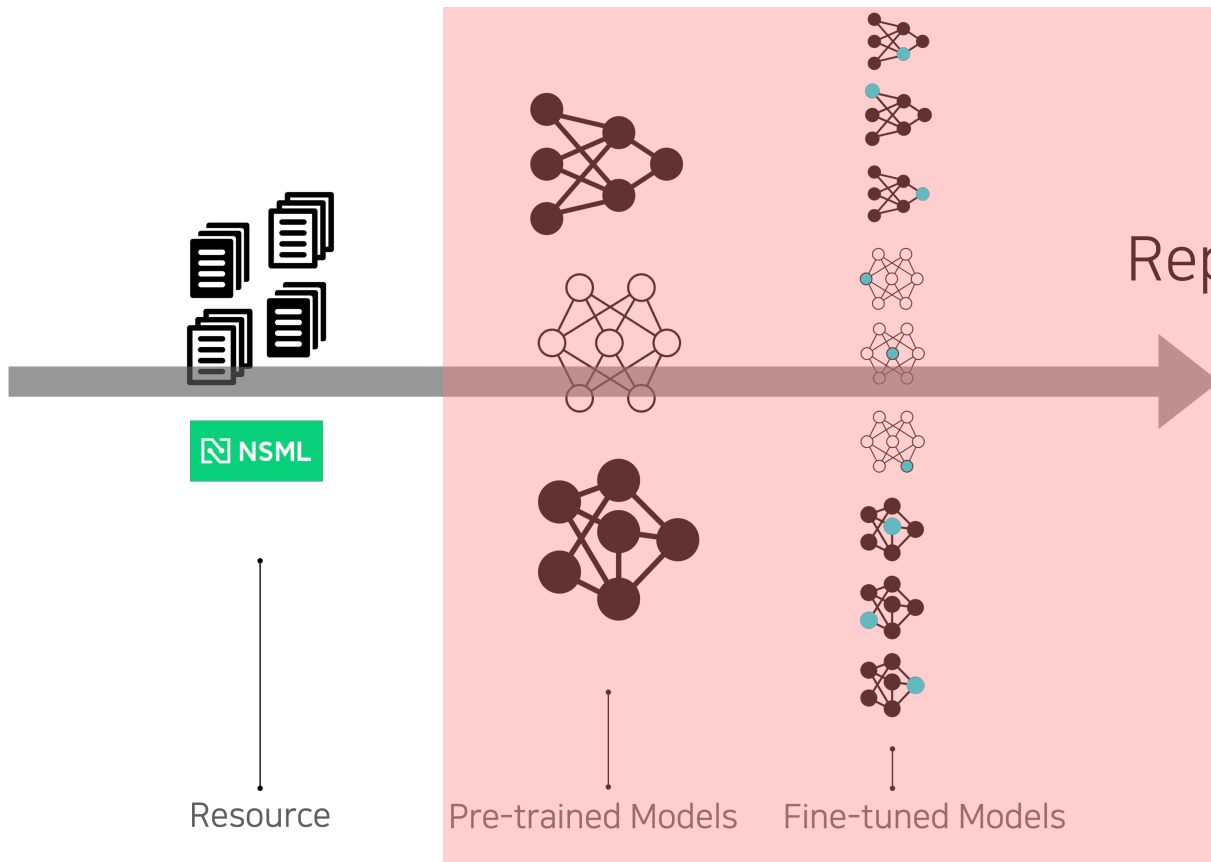


Resource

Pre-trained Models

Fine-tuned Models

# PLM 생산과 사용 라이브러리

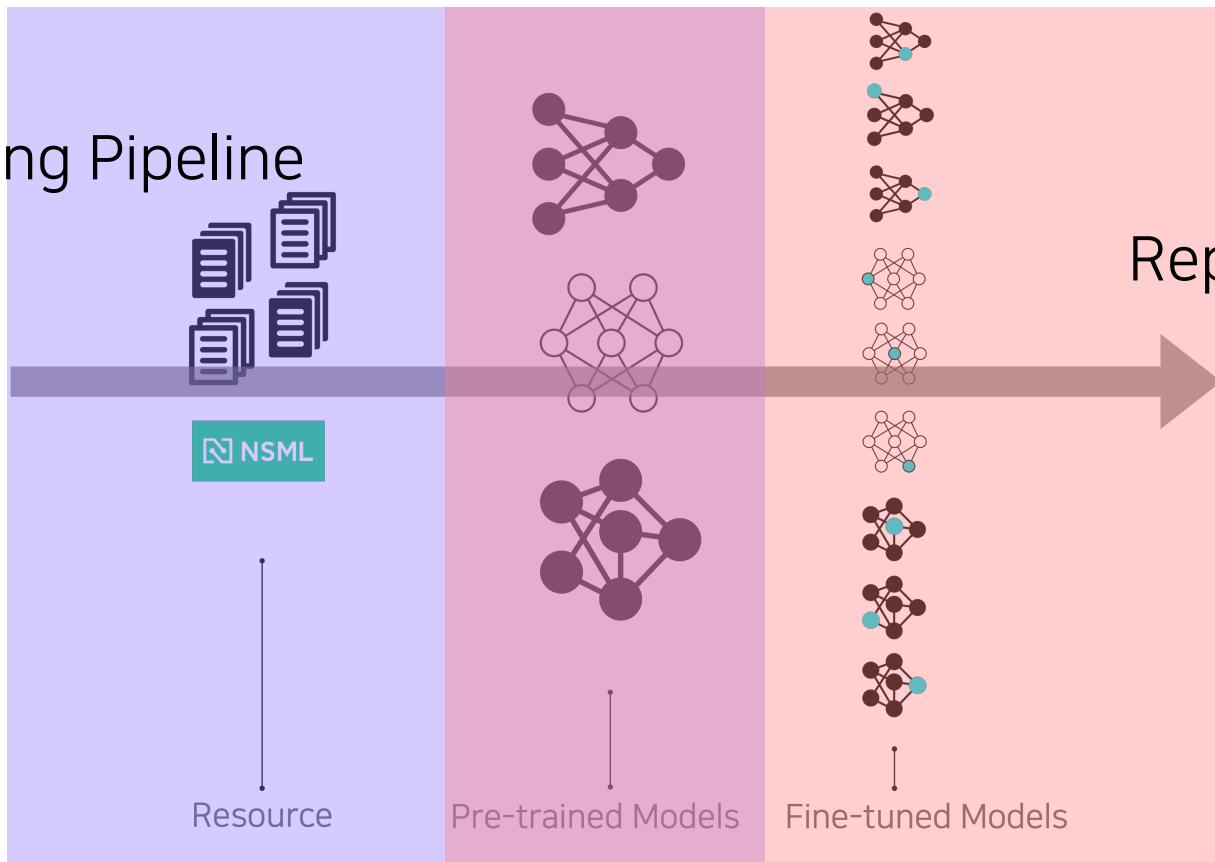


LaRva:  
Language  
Representation  
by Clova



# PLM 생산과 사용 라이브러리

PPL:  
Pre-training Pipeline



LaRva:  
Language  
Representation  
by Clova

# 2.PPL: Pre-training Pipeline

# Pre-train을 하려면...

## 논문에서 공개된 코드 사용하기

- 모델마다 다른 코드 베이스
- 사용 가능한 리소스 제한
- 훈련 완료된 모델 사용자에게도 부담

 [google-research / text-to-text-transfer-transformer](#)

 [google-research / electra](#)

 [pytorch / fairseq](#)

⋮

# Pre-train을 하려면

## 논문에서 공개된 코드 사용하기

- 모델마다 다른 코드 베이스
- 사용 가능한 리소스 제한
- 훈련 완료된 모델 사용자에게도 부담

 [google-research / text-to-text-transfer-transformer](#)

 [google-research / electra](#)

 [pytorch / fairseq](#)

⋮

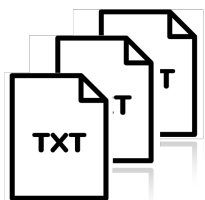


```
python run_pretrain.py \  
pretrain/runner=electra-large \  
data.dataset_path=datasets \  
data.tokenizer_dir=wordpiece_tokenizer
```

# 여러 모델을 쉽게 pretrain 할 수 있도록

## 통일된 인터페이스

- 다양한 모델 등을 조합, 선택 가능하도록 일반화



데이터

WordPiece  
SentencePiece  
ByteLevelBPE  
⋮

토큰라이저

BERT  
ELECTRA  
RoBERTa  
GPT2  
T5  
⋮

모델

Local	Data Parallel	Mixed precision
NSML	Distributed Data Parallel	Gradient accumulation
	DeepSpeed / FairScale	⋮
	Model Parallelism	

훈련

# 여러 모델을 쉽게 pretrain 할 수 있도록

## 통일된 인터페이스

- 다양한 모델 등을 조합, 선택 가능하도록 일반화



```
defaults:
- hydra: pretrain
- pretrain/data: corpus
- pretrain/runner: electra-small
- pretrain/train: ddp
```

```
# @package data
tokenizer_dir: wordpiece_tokenizer
dataset_path: dataset
per_device_batch_size: 32
num_workers: 8
```

```
# @package runner
_target_: ppl.electra.ElectraRunner
model:
  attention_probs_dropout_prob: 0.1
  embedding_size: 1024
  hidden_act: gelu
  ...
  divide_factor: 4
  d_lambda: 50
```

```
optimizer:
  learning_rate: 2e-4
  num_warmup_steps: 10000
  num_training_steps: 400000
  ...
  betas2: 0.999
  weight_decay: 0.01
```

```
# @package _global_
trainer:
  _target_: pytorch_lightning.Trainer
  accelerator: ddp
  gpus: 8
  gradient_clip_val: 1.0
  ...
  precision: 16
  resume_from_checkpoint: null

callback:
- _target_: ppl.utils.SaveModelCheckpoints
  ...
```

# 데이터

## Serialization

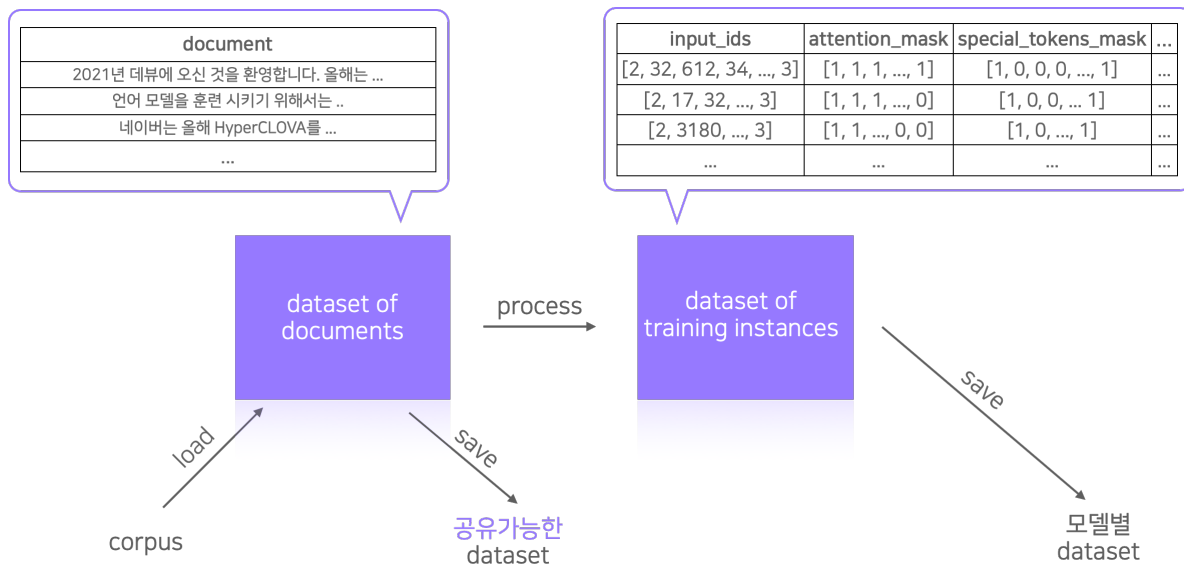
- Raw Corpus를 학습에 사용할 수 있는 형태로 바꾸기
- 훈련 시간 단축



# 데이터

## Pipeline

- 최종 결과는 모델 별로 하나씩
- 데이터 로딩, 전처리 과정은 모델간 공유

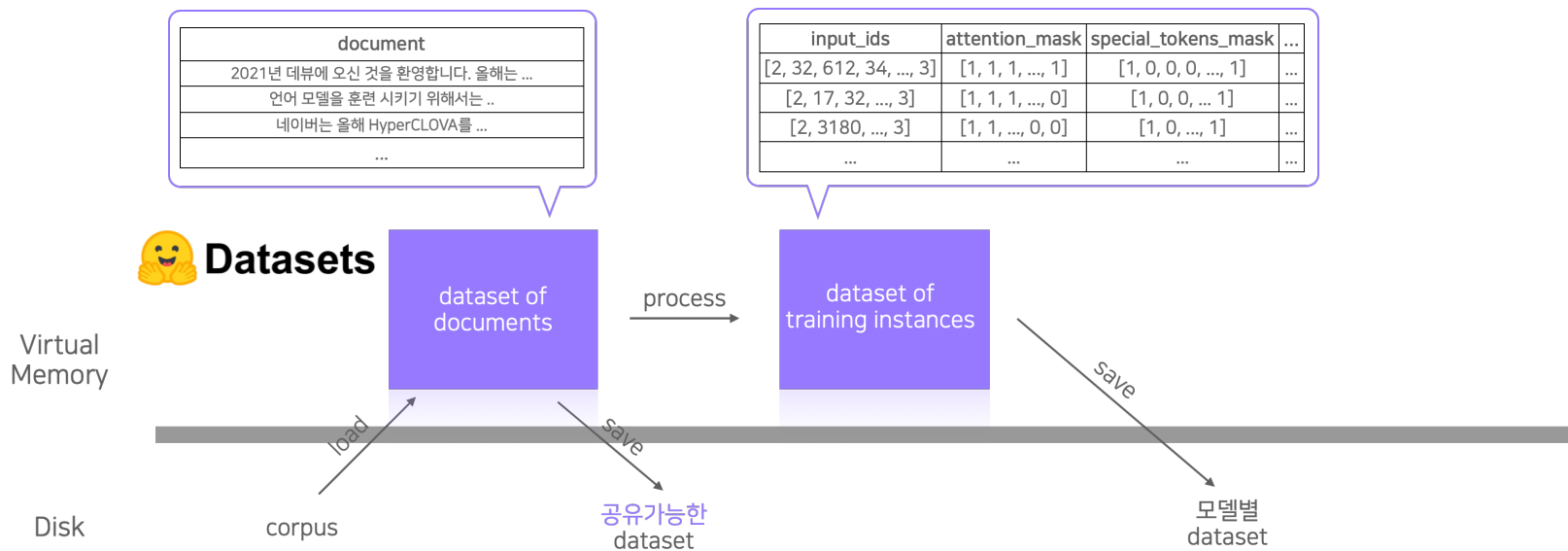




# 데이터

## Arrow 데이터

- 메모리를 (거의) 사용하지 않음
- On memory 데이터처럼 디버그 가능



# 토큰나이저

## Vocabulary 학습 🤗 Tokenizers ❤️

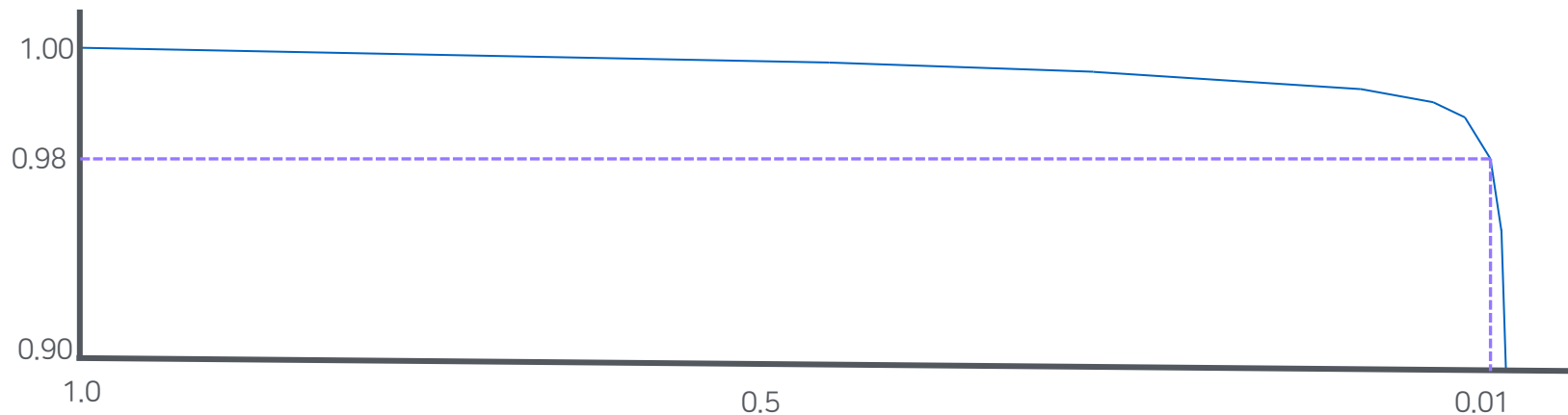
- 데이터 컬렉션 마다
- 형태소 정보가 반영된 BPE 알고리즘



# 토크나이저

## BPE 메모리 사용량 이슈

- 훈련 데이터 샘플링



샘플링 비율에 따른 어휘 중복 비율

# 토큰나이저

## BPE 메모리 사용량 이슈

- 훈련 데이터 샘플링
- 코퍼스 포맷 전처리



Narsil commented on Sep 18, 2020

Collaborator

Hi,

What version of tokenizers are you running ?

BPE algorithm can be quite memory intensive when the length of the tokens is large, which can be the case in Japanese because of no spaces.

We also are making some changes to lower the current memory footprint of some preprocessing.

1. You could try to add some PreTokenizer to split your incoming sentences into relevant groups to try and lower the size of the chunks that get fed to BPE. A simple way would be to preprocess your data and put each split on different lines.
2. In the not so far future, you will be able to train with SentencePiece which notably behaves better on languages that don't have spaces. It's still very bleeding edge.
3. Our memory footprint will be lower in 0.9 so if you can afford to wait that's a solution.
4. If you can't afford to wait and methods above don't work, please tell us, I'm willing to help to get it working with bleeding edge stuff.

<https://github.com/huggingface/tokenizers/issues/422>

## 다양한 모델

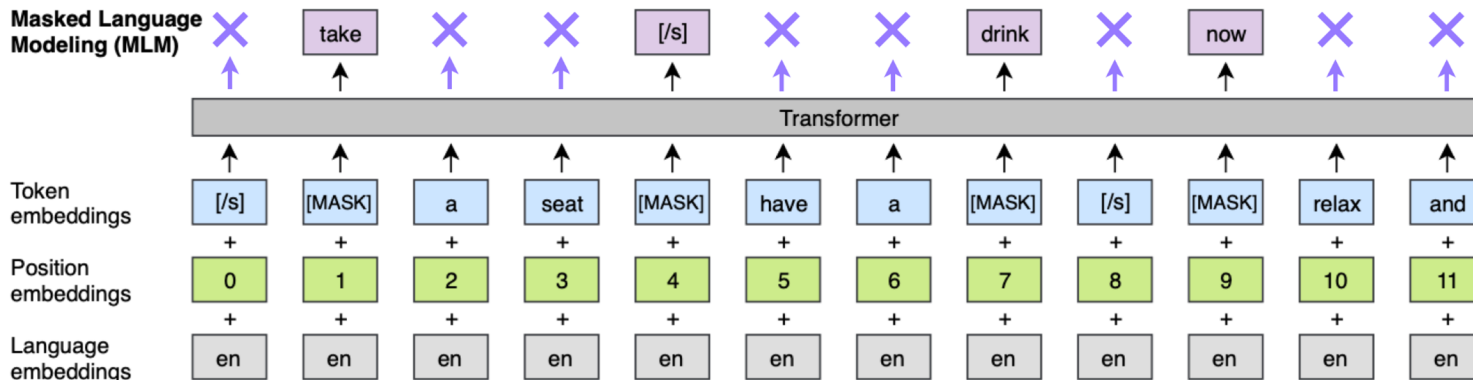
- 🙌 Transformers ❤️




## 불필요한 연산을 줄이기

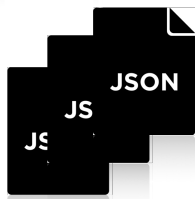
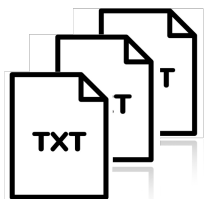
- Loss 계산시 불필요한 부분 제거 가능

✕ : 계산 안 해도 되는 부분



## 더 큰 모델, 더 큰 배치 사이즈를 위해

- DDP with  PyTorch Lightning
- Mixed precision, Gradient accumulation



데이터

WordPiece  
SentencePiece  
ByteLevelBPE  
⋮

토큰라이저

BERT  
ELECTRA  
RoBERTa  
GPT2  
T5  
⋮

모델

Local Data Parallel Mixed precision  
NSML Distributed Data Parallel Gradient accumulation  
DeepSpeed / FairScale ⋮  
Model Parallelism

훈련

# 3.LaRva: Language Representations by Clova



# 모델 학습 이후

모델 학습 이후 어떤 기능이 필요할까?

# 모델 학습 이후

모델 학습 이후 필요한 기능?

PLM  
로드

파인 튜닝

인퍼런스

문장  
임베딩

# 모델 학습 이후

모델 학습 이후 필요한 기능?

PLM  
로드

파인 튜닝

인퍼런스

문장  
임베딩

# 모델 학습 이후

모델 학습 이후 필요한 기능?

PLM  
로드

파인 튜닝

인퍼런스

문장  
임베딩

# 모델 학습 이후

## 무엇을 고려해야 할까?

- PLM 및 사내 데이터셋 배포: 기존 통일된 방법 부재
- 확장성있는 튜닝 예시
  - 데이터에 따라 전처리 등 모델러가 직접 코드 개발
  - HyperCLOVA 이후 새로운 방법론

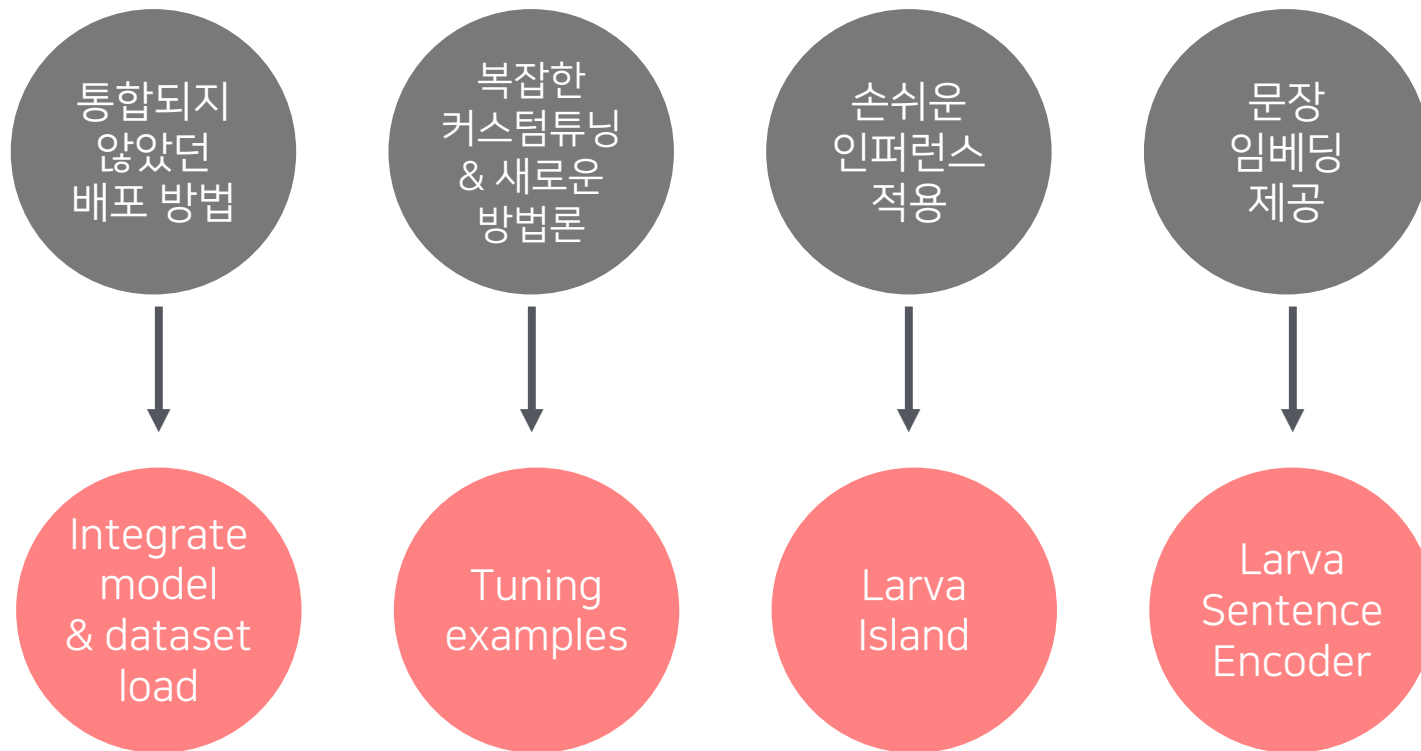


# 모델 학습 이후

## 무엇을 고려해야 할까?

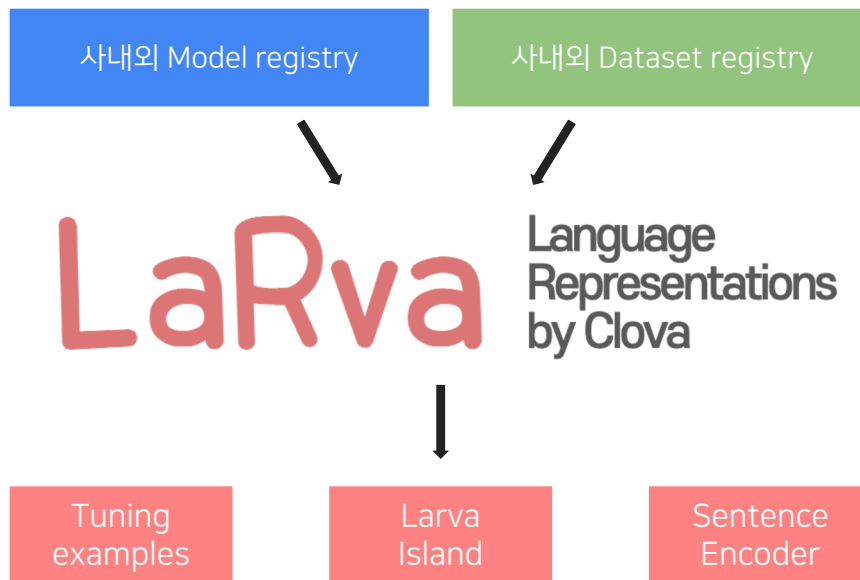
- 손쉬운 인퍼런스
  - 인퍼런스만 손쉽게 하고 싶은 사용자를 위한 프레임워크
- 문장 임베딩 제공
  - 다양한 모델 학습 방법 존재
  - 서비스에 적용 가능하도록 다양한 기능 제공

# LaRva 라이브러리



# LaRva 라이브러리

사내외 PLM 및 데이터셋 로딩 방식 통합

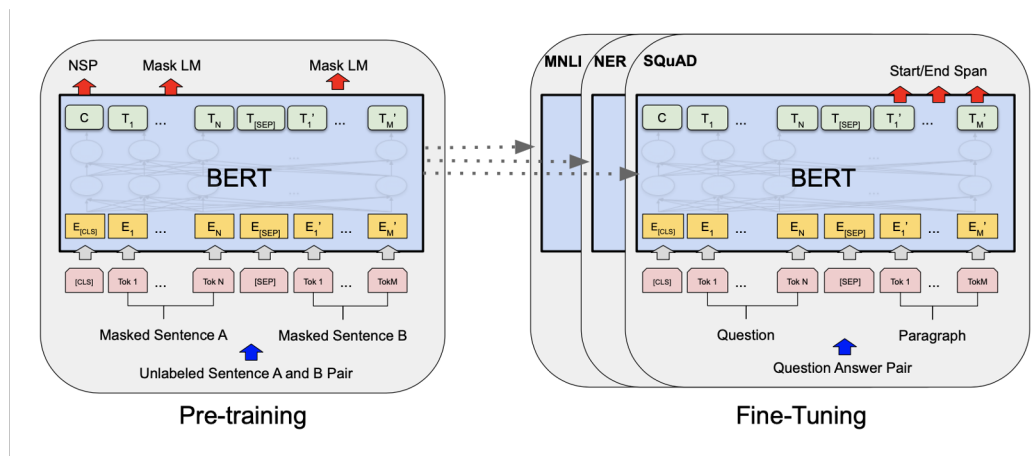




# Tuning before HyperCLOVA

## HyperCLOVA 이전의 모델을 활용한 fine-tuning

- Pre-training 단계에서 배운 지식의 "transfer"
- 모델의 파라미터 미세 조정



# Tuning before HyperCLOVA

## 확장성있는 fine-tuning 예시

- 데이터마다 서로 다른 포맷 & 전처리 🥲
- Base 클래스 통해 기본적인 Preprocess, Metric 등 정의
- 최소한의 코드만 모델러가 수정하도록 확장성 고려



# Tuning after HyperCLOVA

## 모델 사이즈 증가에 따른 fine-tuning 방법론의 한계점

- Large scale 모델의 전체 파라미터 튜닝하기엔 어려움
- 새로운 방법론 P-tuning, LoRA를 tuning examples로 제공

In-  
context  
Learning

P-tuning

LoRA

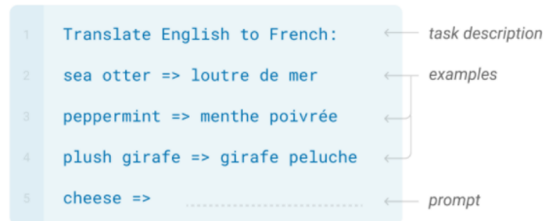
# Tuning after HyperCLOVA

## In-context Learning

- Few-shot
- Prompt engineering

### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

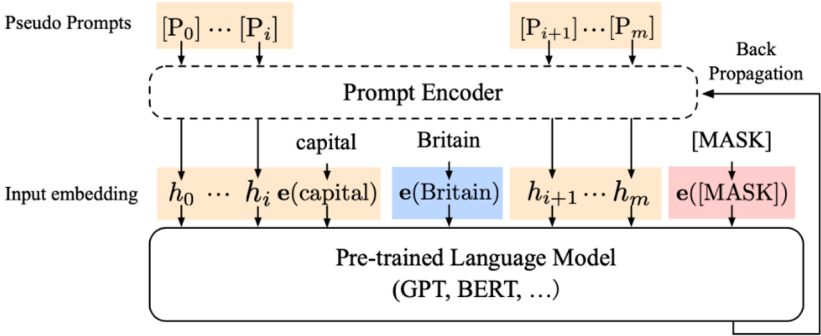
### Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



# Tuning after HyperCLOVA

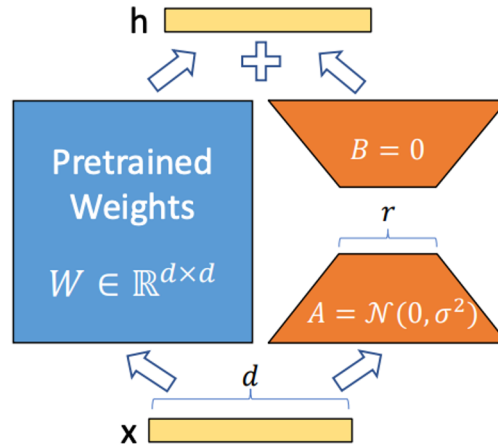
## Prompt based tuning (P-tuning)



(b) P-tuning

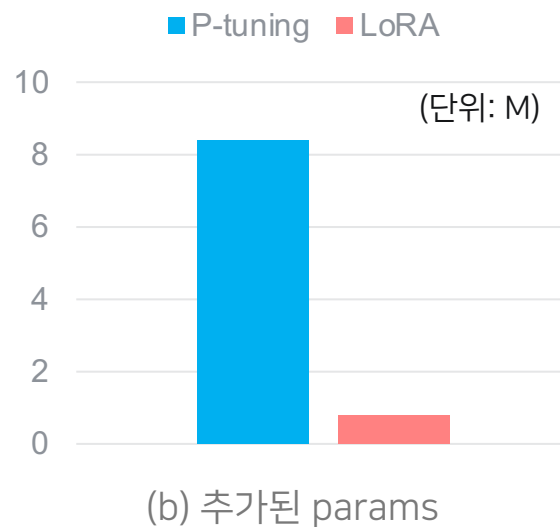
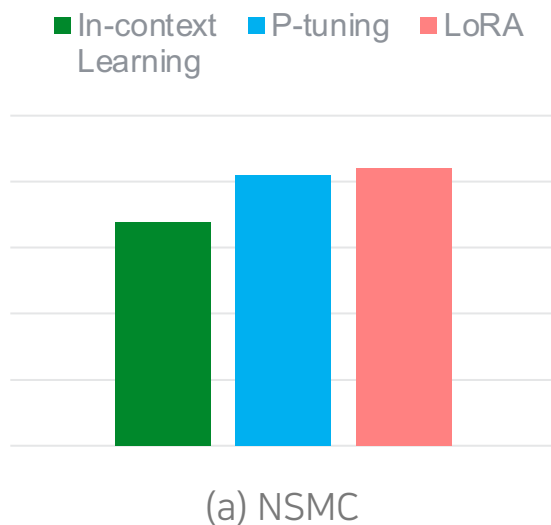
# Tuning after HyperCLOVA

## LoRA: Low-Rank Adaption



# Tuning after HyperCLOVA

## HyperCLOVA 1.3B 실험 결과



# Larva Island

## 튜닝 이후 인퍼런스

- Larva Island: 손쉬운 인퍼런스를 위한 프레임워크
- 튜닝 모델 변경만으로 인퍼런스 가능: 손쉬운 연구/서비스 적용
- 웹데모 및 API 사용 가능
- 배치 단위 처리 가능

```
from larva import LarvaIsland  
  
island = LarvaIsland("sentiment-analysis", model="model_code")  
  
island("오늘 데뷰 발표 모두 재미있었어요!")
```

다음 빈 칸에 문장을 입력하면 해당 문장의 개체명(Named Entity)을 인식합니다.

네이버 본사는 경기 성남시 분당구에 있다.

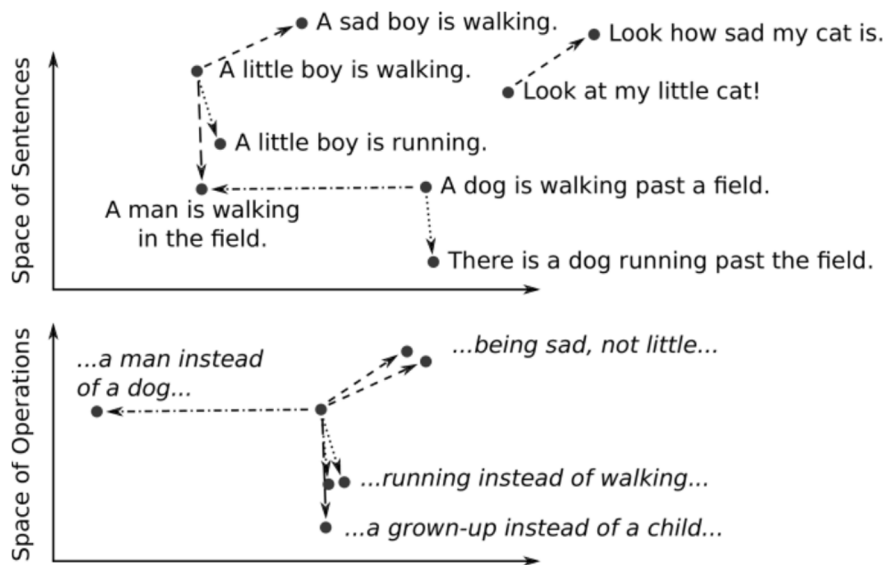
네이버 **ORG** 본사는 경기 성남시 분당구에 **LOC** 있다.



# Sentence Encoder

## 문장 임베딩의 필요성: 다양한 사용 가능성

- Retrieval
- Clustering
- Classification



# Sentence Encoder

## 사용성을 위한 인터페이스

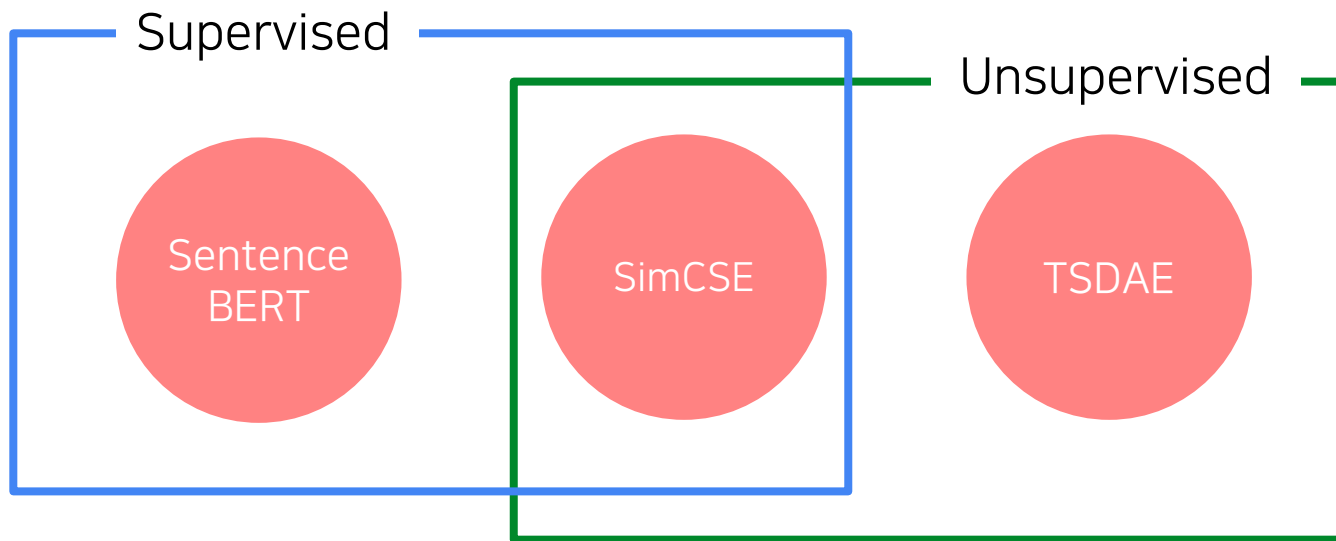
- 손쉬운 사용을 위해 허깅페이스 🤗 구조와 동일하게 맞춤
- 문장 임베딩 생성, Similarity search 등의 기능 제공
- Faiss, 리소스 선택 등의 옵션 추가 제공

```
from larva import LarvaSentenceEncoder

encoder = LarvaSentenceEncoder.from_pretrained("model_code", use_faiss=False)

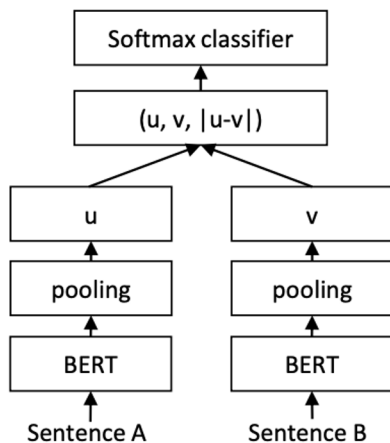
# generate embedding
encoder.encode(["강아지가 산책을 하고 있다.", "고양이 한 마리가 낮잠을 자고 있다."])
```

# Sentence Encoder

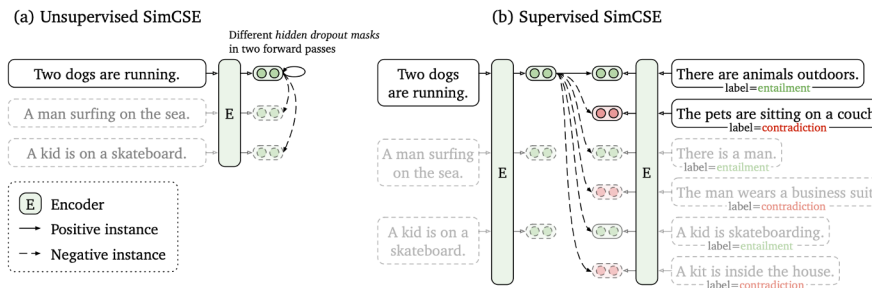


# Sentence Encoder

## Sentence-BERT



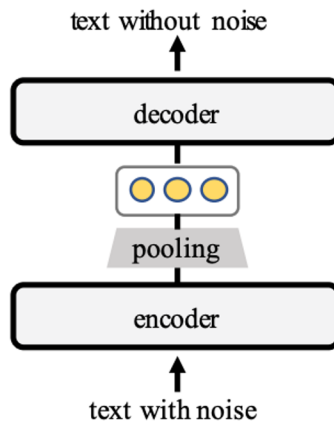
## SimCSE



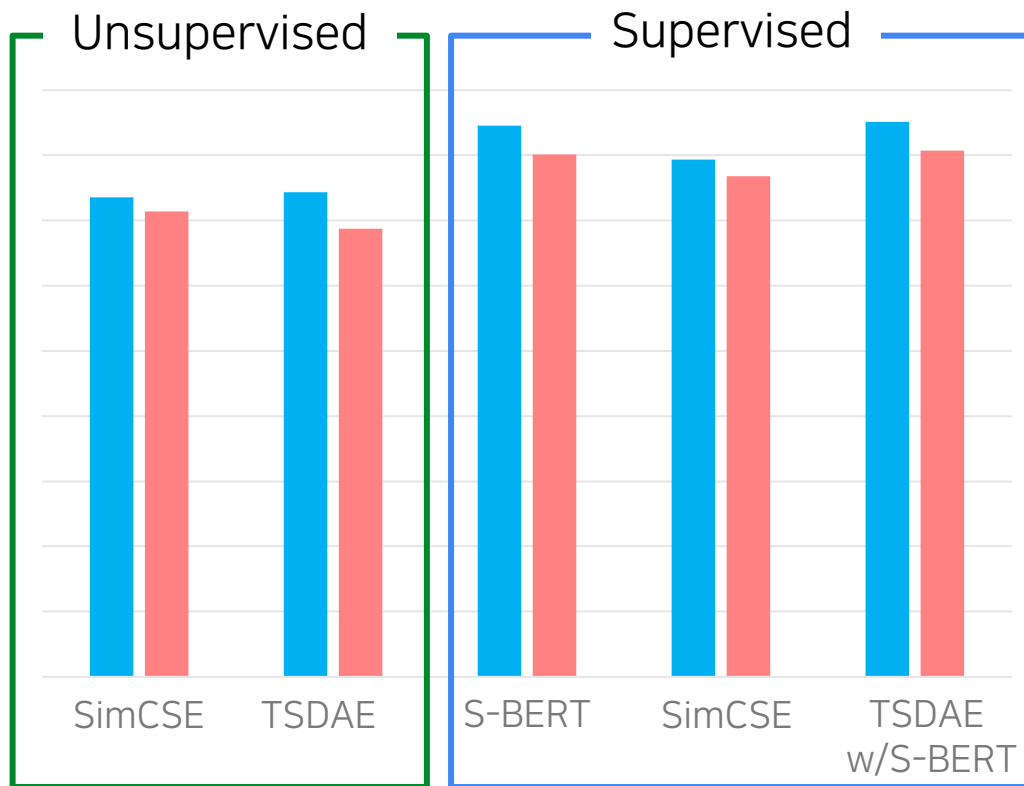
# Sentence Encoder

## TSDAE

- Transformer-based Sequential Denoising Auto-Encoder
- 오토 인코더 형태로 노이즈 문장을 복원하도록 학습



# Sentence Encoder

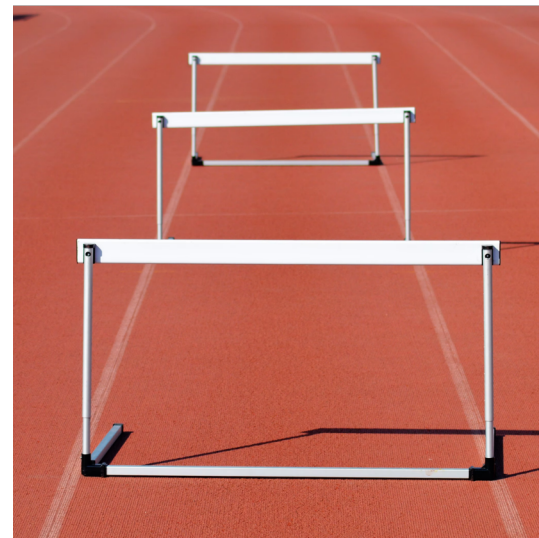


# 4. Future Work

# Challenges

## 앞으로 넘어야 할 문제들

- 실제 서비스 적용? 큰 사이즈의 모델은 아직 어려움이 있음
  - 모델 Parallelism
  - 모델 경량화 - Distill, Quantization



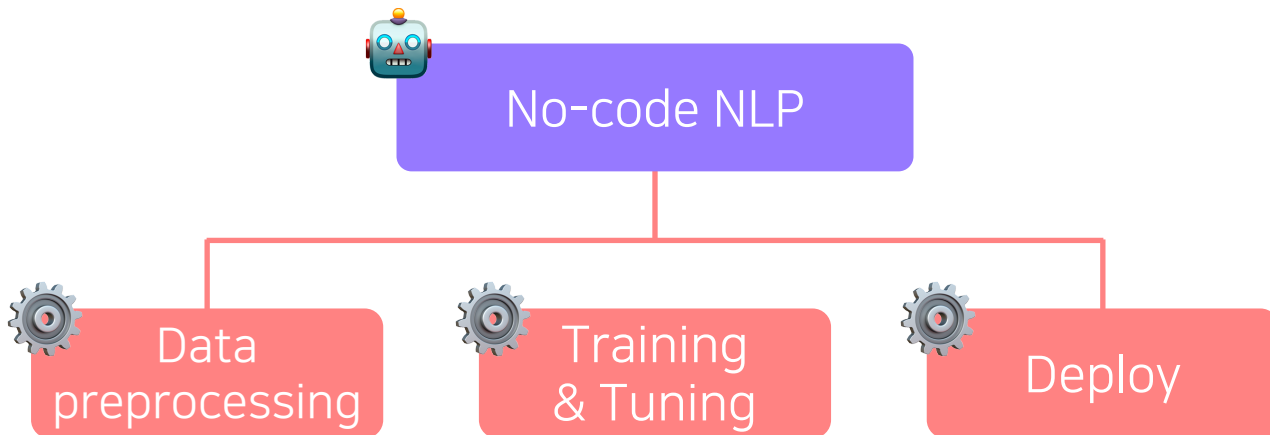


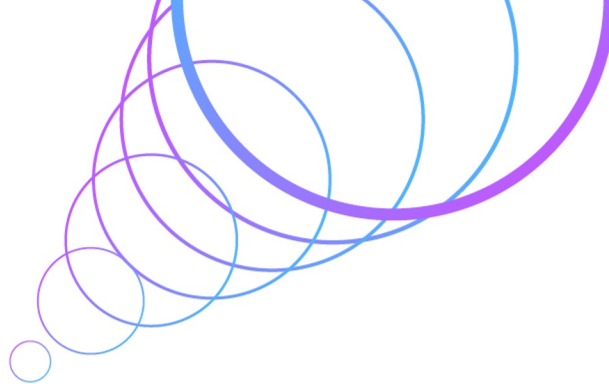


# No-code NLP

## 손쉬운 No-code NLP

- 데이터 전처리, 학습, 배포 등의 자동화: No-code
- 작은 스케일의 LM부터 Large scale의 LM까지 포괄





Thank You

