# Parameter Efficient Transfer-Learning for NLP. \<paper\>

## \<abstract\>

fine-tuning은 다른 다운스트링 테스크가 들어올때마다 기존 parameter/model 재학습시켜야 함음. 이 대신 *adapter module*을 활용한 transfer을 소개함.

기존의 학습된 모델의 parameter을 그대로 두어 공유하고, 새로운 레이어를 추가함.

## \< Intro \>

BERT: large text corpora로 마스크 loss와함께 학습된 Transformer Network. (Text classification/ 회귀(Qk))

목표는 BERT로 다양한 task를 수행할 수 있는 전이학습 application을 만드는것 → Compact / Extensible downstream Model

① Compact : 각 테스크에 적은 라메메터의 추가만으로 많은 테스크 해결가능.
② Extensible : 전의 것을 잊지 않고 새로운 테스크 해결가능.

이전의 전이학습은 fine-tuning과 feature-based transfer learning으로 나눔. 여기서 adapter module을 새로 제안함.

↳ 사전학습된 real-valued embedding vector를 사용.
↳ 사전학습 network에서 weight를 copy, downstream에 맞게 튜닝.
⇒ 각 테스크마다 새로운 weight set이 필요함. 하지만 fine-tuning보다 adapter은 두배가까이 적은 param으로 비슷한 성능 보임.

Adapter? : 사전 학습된 네트워크의 층사이에 추가되는 새로운 모듈. *base model + adapter module*

$\phi_W(x)$  function with parameters $W$
- FB : $\chi_v(\phi_W(x))$ - 새로운 함수 $\chi$, $V$만 학습.
- FT : compactness를 서향하여 각 테스크에 $W$을 조정.
- AT : $\psi_{W,V}(x)$ , 새로운 함수 정의, $W$는 copy, 초기 파라미터 $V_0$는 $\psi_{W,V_0}(x) \approx \phi_W(x)$에서 온. $V$만이 학습.

↳ $W$는 쟁긴이어의 테스트에 전의 모델에 영향 주지 X. 보통 기존 network $\phi_W$에 새로운 레이어를 추가하는 방식. (param)

① 수식적차이 against Feature-Based/ Fine-tuning.

② 방법의차이 against Multitask/Continual
- Multi-task Learning : 모든 task가 동시적인 접근 요구 ⇒ AT는 아님.
- Continual Learning : 끊임없는 task 처리, 하지만 재학습을 줄이는 전의 task를 잊어버림.

⇒ 반면 AT는 공유된 parameter에 대한 동결하기때문에 완벽하게 전 task를 기억함.

## \<Adapter tuning for NLP\>

Goal : ① 좋은 성능. ② 데이터에 순차적접근이 아닌 순차적 학습 ③ 각 테스크에 적은 수의 파라미터만 추가하는 것

Deep network의 vanila Fine-tuning에서는 네트워크의 가장 위 layer에 수정이 가해짐. 이는 upstream과 downstream task의 label space와 loss가 다르기 때문에 필요반영하게 됨.
하지만 Adapter tuning 에서는 가장 위의 layer를 하나 더추가하고 이 layer의 parameter를 random으로 초기화시켜 학습시킴.
이는 테스크가 추가되어도 총 model size의 증가율을 줄여주고, Near-Identity initialization을 Adapter에 적용하여 기존의 network에 영향을 주지 않고 좋은 성능의 학습을 유도.

↳ layer 하나씩만 추가. ↳ 그 layer에 가까운 값으로 param 초기화. 기울기 소실, 폭발 문제를 예방.

## \<Instantiation for Transformer Networks\>



- Feed forward와 Attention layer 모두 projection을 통해 input output size로 맞춤.
- Both use skip-connection. ( at Fc레 / 일그니어 가까워야할 차이를 해결) → Normalization
- 각 sub-layer 뒤에 adapter를 직접적으로 apply함 (projection 후, skip-connection 전), 그 뒤 바로 Normalization에 입력으로.

- parameter 수를 제한하기 위해 *bottleneck architecture* MLP. (특징 압축, 차원축소)
  o d-dim의 feature를 더적은 크기의 m-dim으로 project. (m: bottleneck dimension)
  o 각 레이어는 bias를 포함하여 2md + d+m 만큼의 parameter가 추가됨.
  o 대략으로 기존 모델 파라미터의 0.5% ~ 8% 가량을 사용.
  o Adapter는 skip-connection을 이미 포함. 🔁
  o projection layer의 파라미터가 0에 가깝게 초기화되면, module은 거의 Identity function으로 초기화됨.

- 각 task마다 new layer normalization parameter를 학습. (conditional batch normalization과 비슷).
  o 각 레이어 마다 2d 만큼의 parameter 가 추가됨.

## \<Experiments\>

GLUE를 포함, 다양한 데이터에 적용
- Adam 최적화
- Pre-trained BERT Transformer network. (24 layer, 330M param)

↳ ex) FT: 9x, AT: 1.3x

→ classification task에 있어 최고 성능은 fine-tuning 이지만 적은 parameter만으로 적은 성능차이를 이뤄냄.