

Adversarial Multi-task Learning for Text Classification

Multi-task Learning

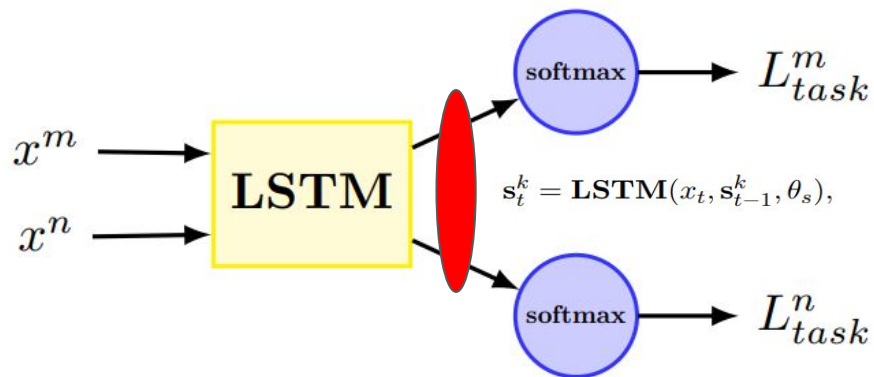
- 여러 Task를 두고 하나의 모델로 학습시키는 학습 방법
- 보통 **overfitting**을 방지하기 위해 **regularization parameter**를 추가하는데, 이보다 자연적으로 다양한 **task**를 수행하기 위해 추가된 **regularization**이 더 효과적일 때가 있음. 이때 **multi-task learning**이 매우 효과적임
- Input - Label pair가 부족한 NLP에서 많이 사용됨

전이 학습과의 비교

- 전이학습은 최종적 **task**를 위해 그 전의 **task**들을 학습 시킴 (보통 한가지 **task**만 한다고 함)
 $t(1), t(2), \dots, t(n-1) \rightarrow t(n)$
- Multi-task Learning은 포함된 모든 **task**를 동시에 학습 시킴
 $t(1), t(2), \dots, t(n-1), t(n)$
- Multi-task Learning을 전이학습의 범주에 포함시키는 경우도 있음
Parallel Transfer Learning : Multi-task Learning
Sequential Transfer Learning : Transfer Learning
- Multi-task Learning으로 학습된 모델도 전이 학습 모델로 사용될 수 있음

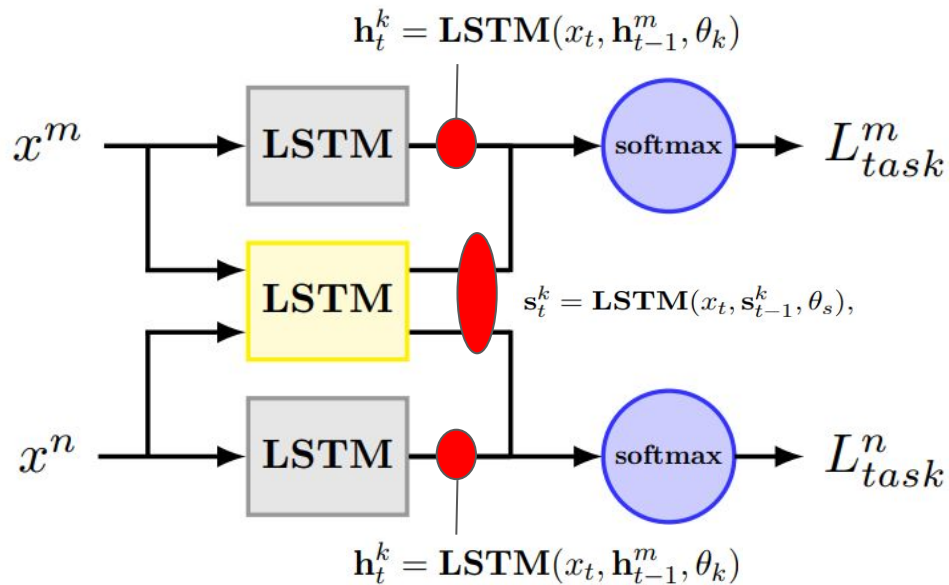
Multi-task Learning

1) Fully-Shared MTL



Multi-task Learning

2) Shared-Private MTL



Objective Function

$$L = L_{Task} + \lambda L_{Adv} + \gamma L_{Diff}$$

L_{Task}

Task-specific softmax

각 Task에 해당하는 Output Layer의 출력 값으로 Loss 값 계산

$$\mathbf{h}_t^k = \mathbf{LSTM}(x_t, \mathbf{h}_{t-1}^m, \theta_k)$$

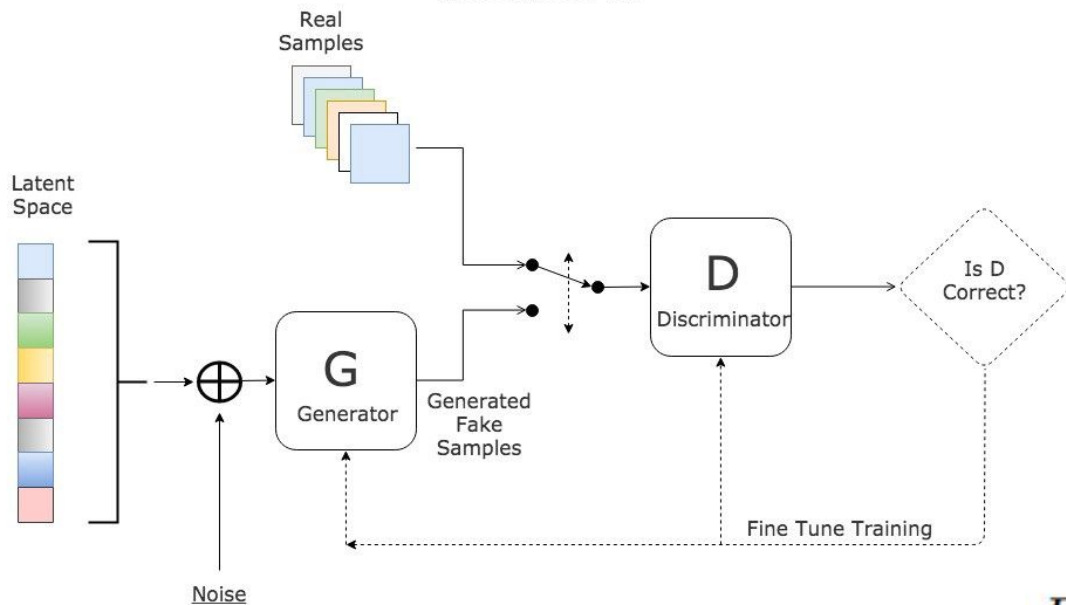
Objective Function

$$L(\hat{y}, y) = - \sum_{i=1}^N \sum_{j=1}^C y_i^j \log(\hat{y}_i^j),$$

$$L_{Task} = \sum_{k=1}^K \alpha_k L(\hat{y}^{(k)}, y^{(k)})$$

$$\lambda L_{Adv}$$

Generative Adversarial Network



Real과 Fake 구분을 못하게 하는 것이 목적

Domain adaptation:
구분이 되지 않으면
transferable feature
= task-invariant feature

$$D(\mathbf{s}_T^k, \theta_D) = \text{softmax}(\mathbf{b} + \mathbf{U}\mathbf{s}_T^k)$$

$$\lambda L_{Adv}$$

Adversarial Loss

Discriminator 모델이 구분을 못하게 Output 출력

$$\mathbf{s}_t^k = \mathbf{LSTM}(x_t, \mathbf{s}_{t-1}^k, \theta_s),$$

Objective function

$$\begin{aligned} \phi = \min_G \max_D & \left(E_{x \sim P_{data}} [\log D(x)] \right. \\ & \left. + E_{z \sim p(z)} [\log(1 - D(G(z)))] \right) \end{aligned}$$

$$L_{Adv} = \min_{\theta_s} \left(\lambda \max_{\theta_D} \left(\sum_{k=1}^K \sum_{i=1}^{N_k} d_i^k \log[D(E(\mathbf{x}^k))] \right) \right)$$

γL_{Diff}

Orthogonality Restriction

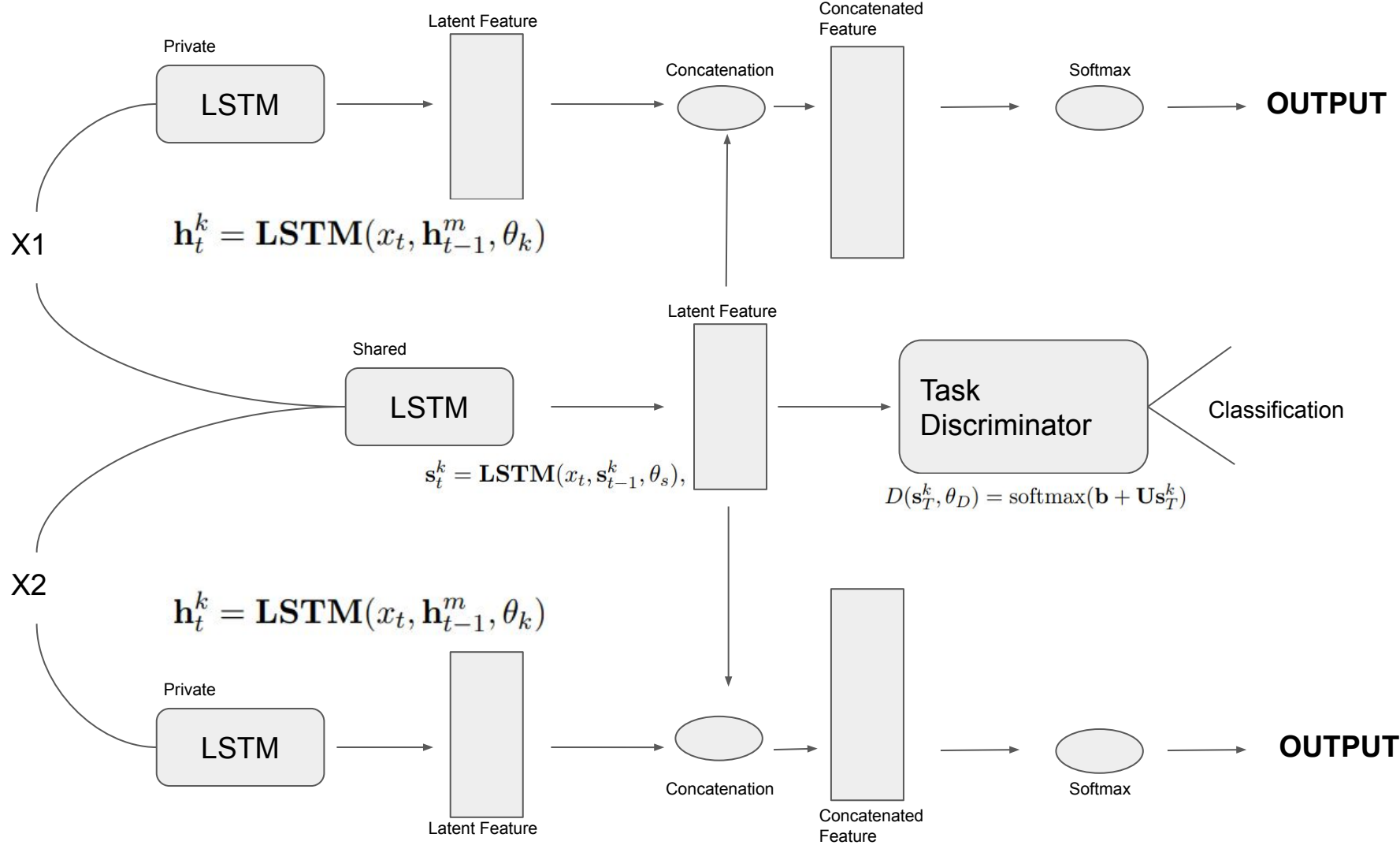
상호배타적 feature representation을 위해 추가 loss 함수를 더함

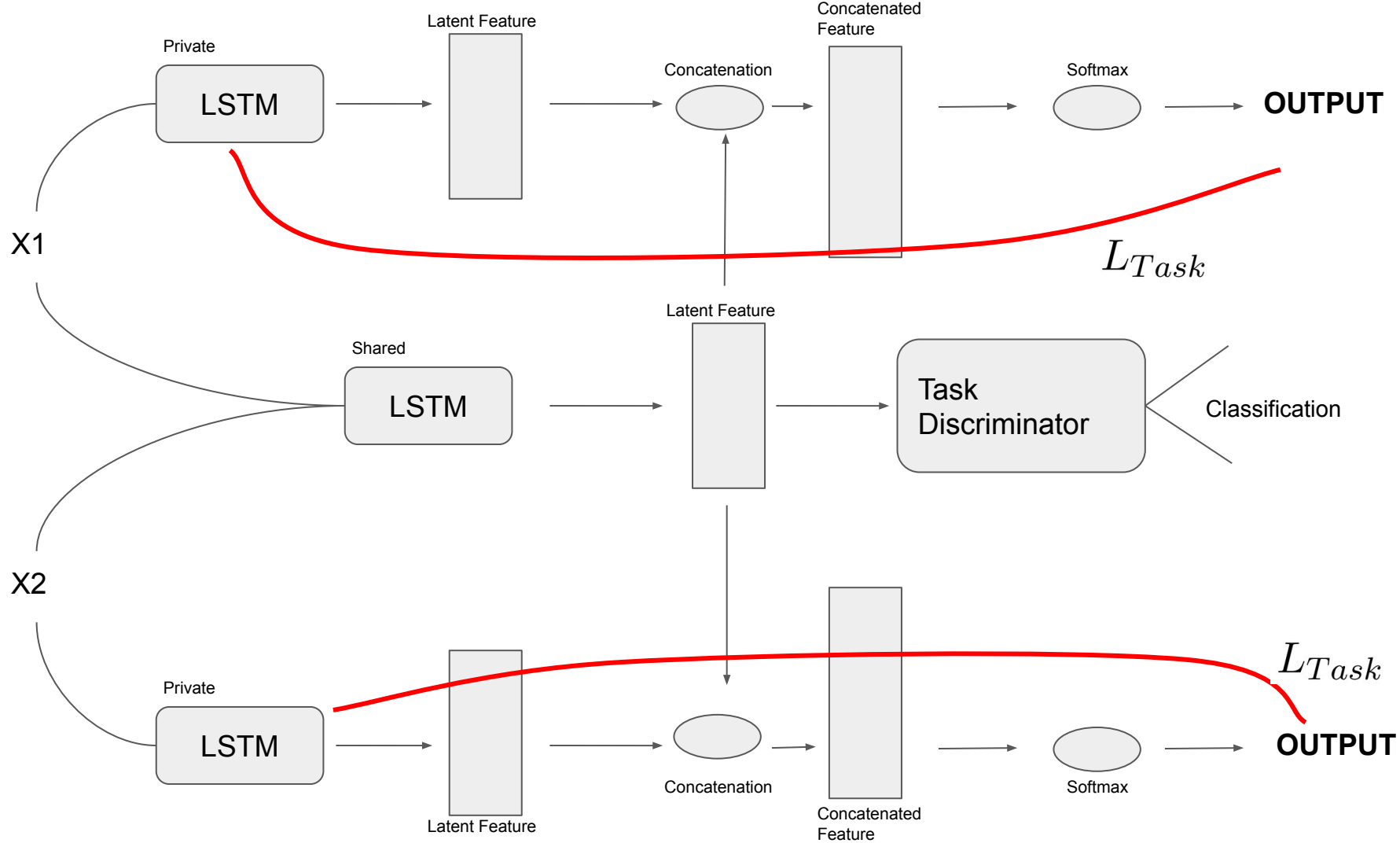
중복 latent representation을 낮춤

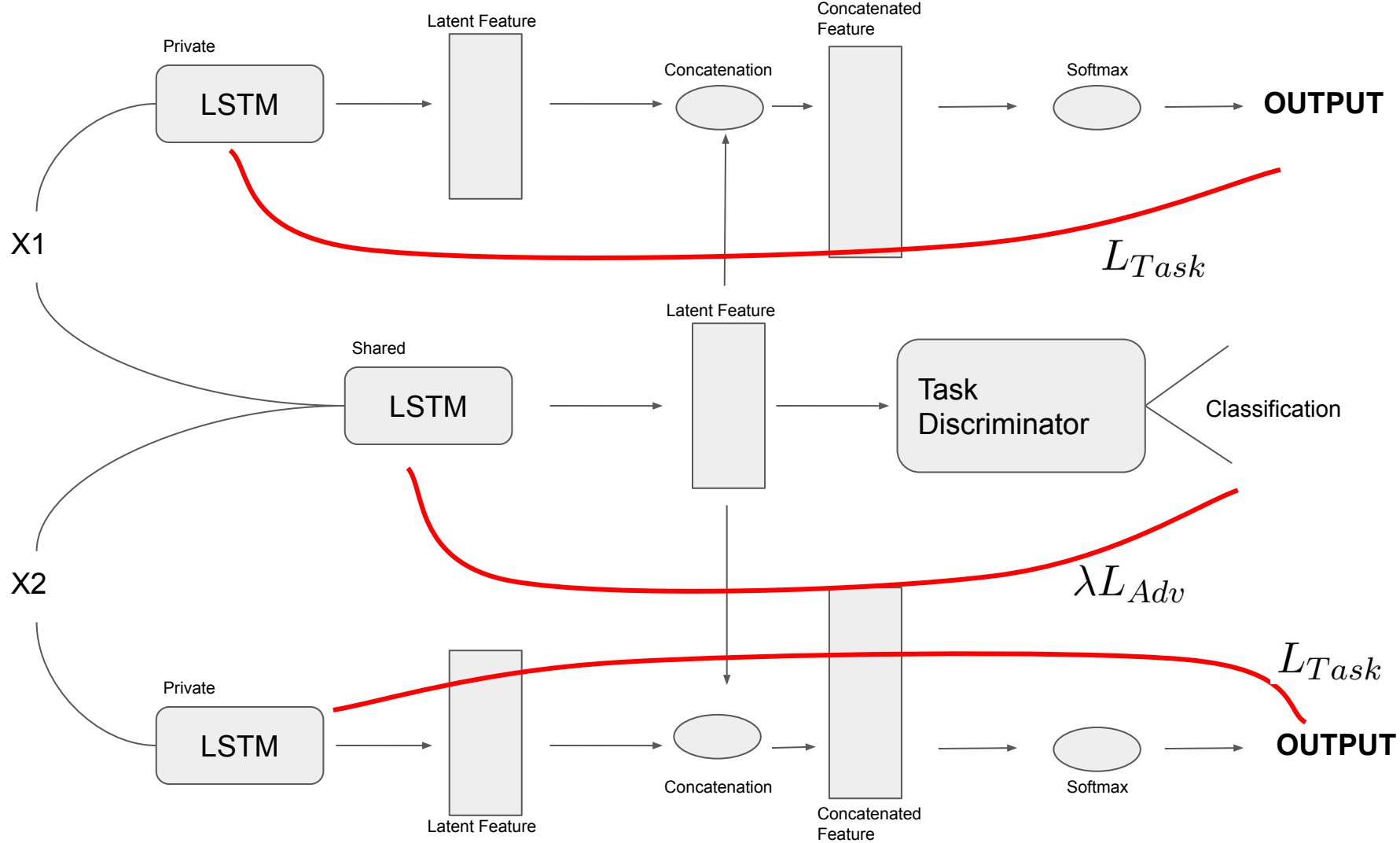
Objective function

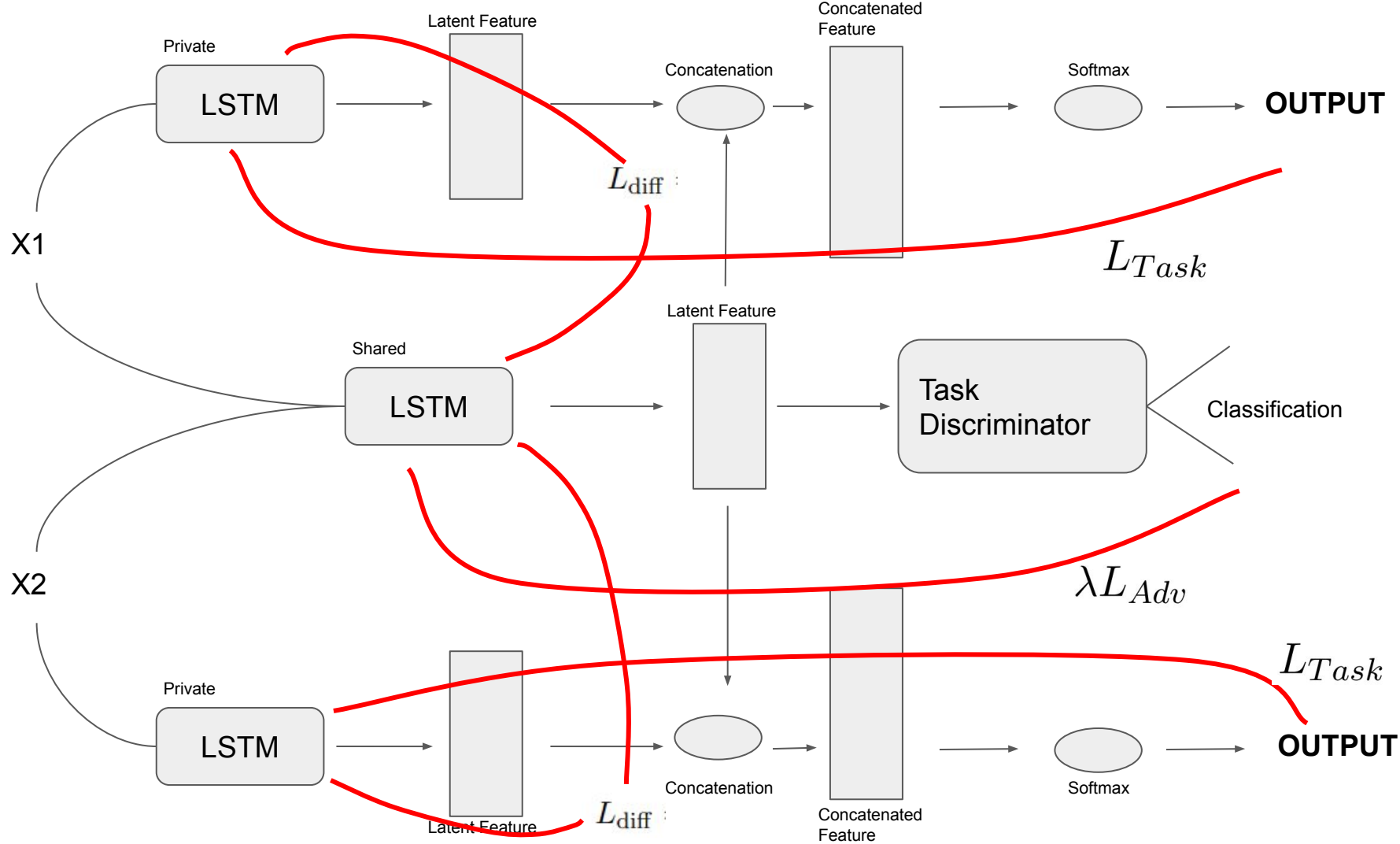
$$L_{\text{diff}} = \sum_{k=1}^K \left\| \mathbf{S}^k{}^\top \mathbf{H}^k \right\|_F^2 ,$$

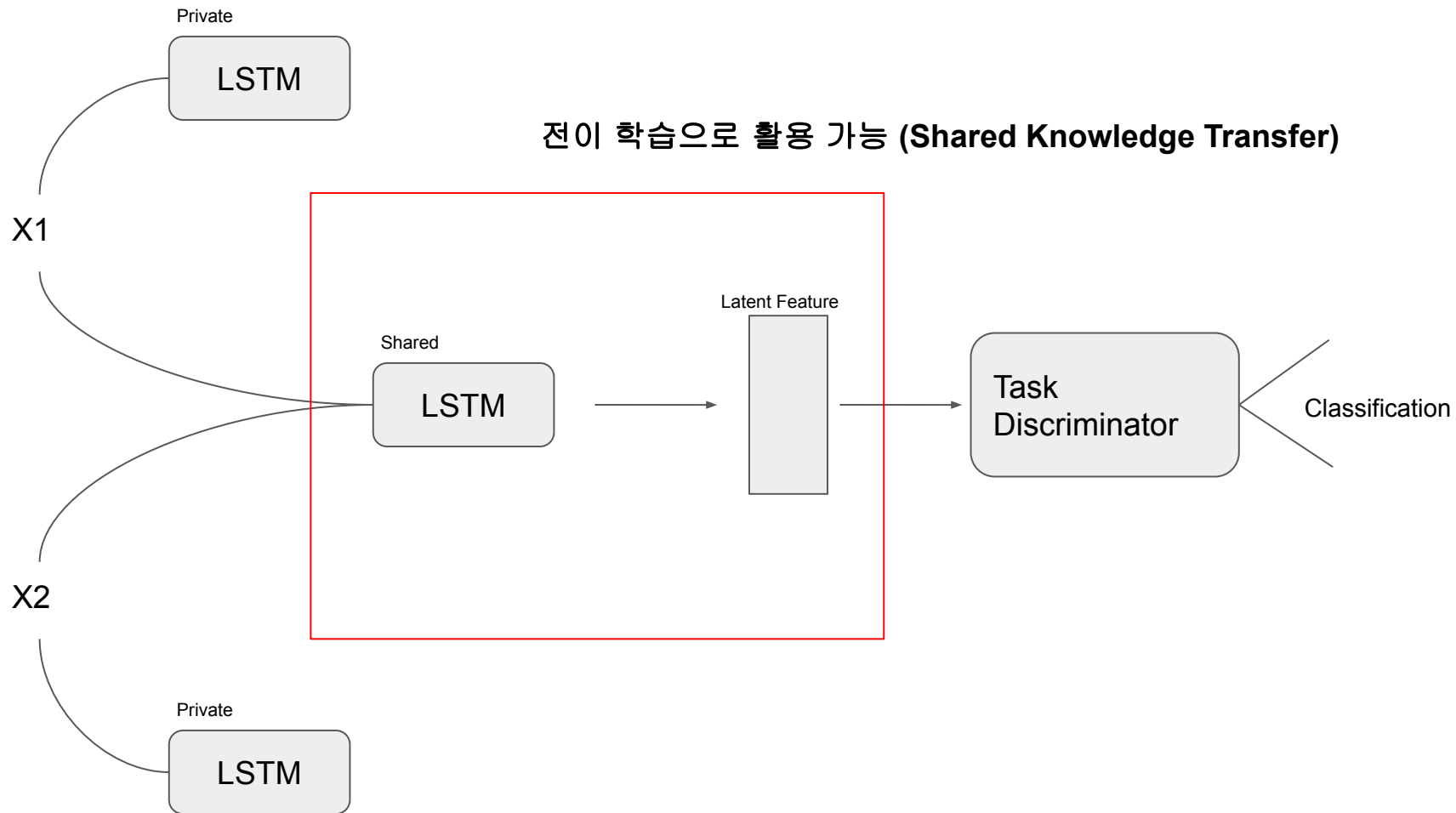
where $\| \cdot \|_F^2$ is the squared Frobenius norm. \mathbf{S}^k and \mathbf{H}^k are two matrices, whose rows are the output of shared extractor $E_s(\cdot; \theta_s)$ and task-specific extractor $E_k(\cdot; \theta_k)$ of a input sentence.



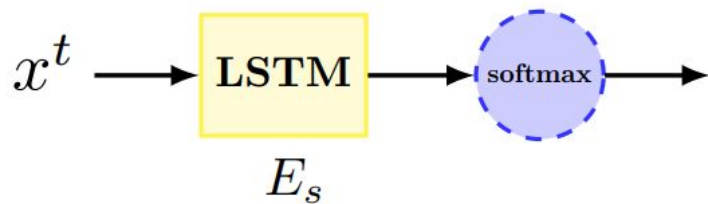




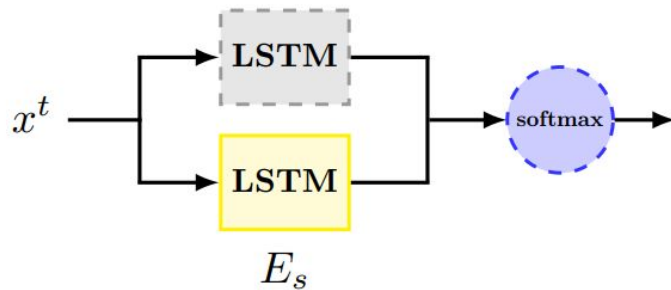




전이 학습 (Shared Knowledge Transfer)

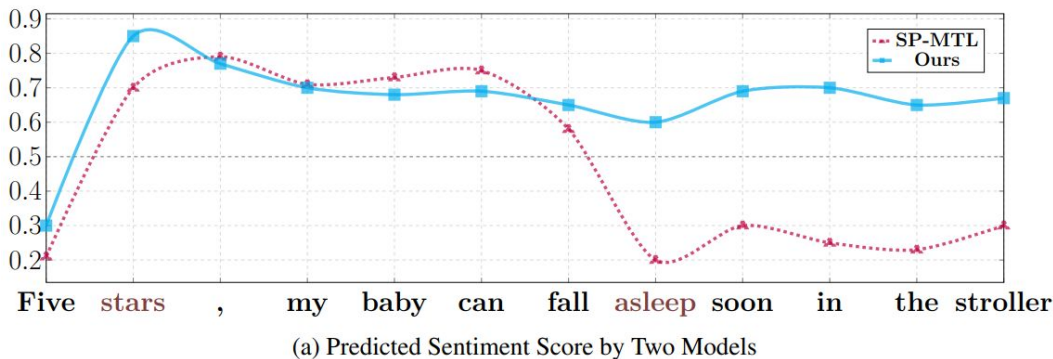


(a) Single Channel



(b) Bi-Channel

Shared FS 와 Private FS의 구분의 중요성



“Asleep”이라는 단어에서 갑자기 부정적인 sentiment로 예측

“Asleep”이라는 단어가 shared feature space에 존재 하고, 다른 task에서는 부정적인 의미로 사용되었다고 추측

사용된 개념

- Adversarial Network (from GAN) <https://arxiv.org/abs/1406.2661>
- Domain Adaptation https://www.alexkulesza.com/pubs/adapt_mlj10.pdf
- Orthogonality Constraints <https://mitpress.mit.edu/books/advances-neural-information-processing-systems>

실제 코드

- <https://github.com/ChenglongChen/tensorflow-ASP-MTL/blob/master/core/model.py>
- https://github.com/FrankWork/fudan_mtl_reviews