Shakil Ibne Ahsan *, Djamel Djenouri and Rakibul Haider
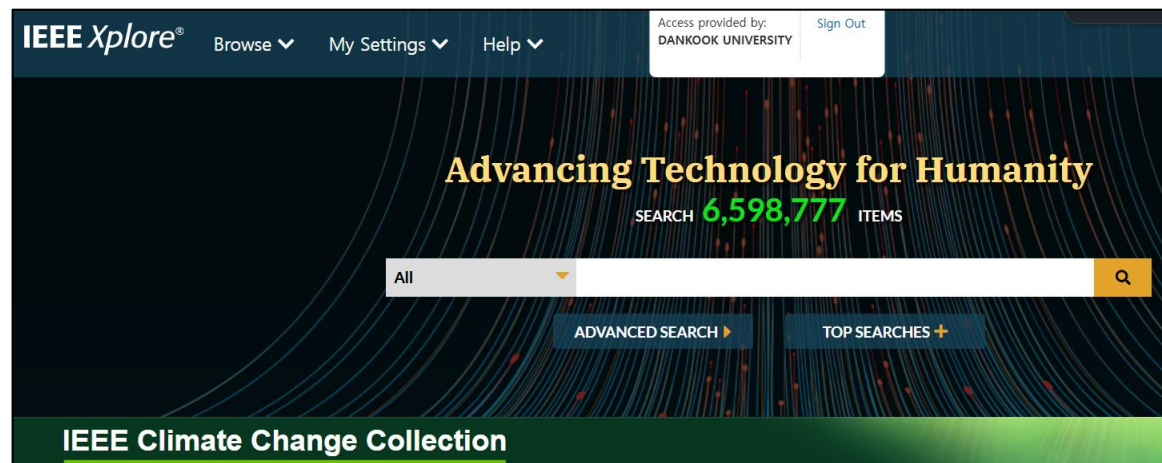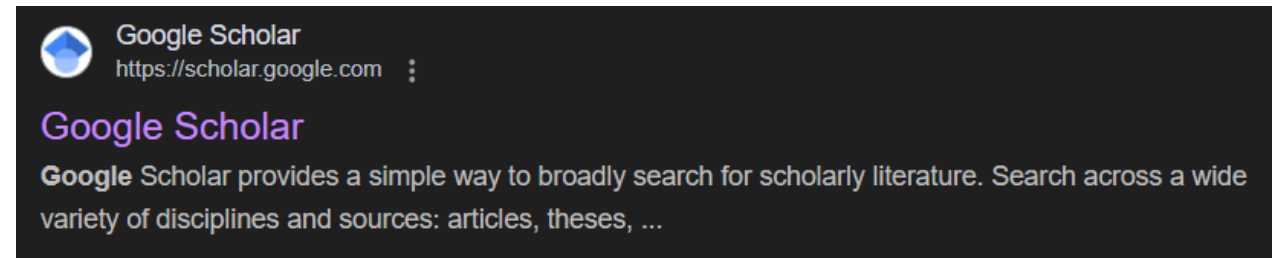
# Federated Learnig with Data Obfuscation and Bidirectional Encoder Representations from Transformers Paper review
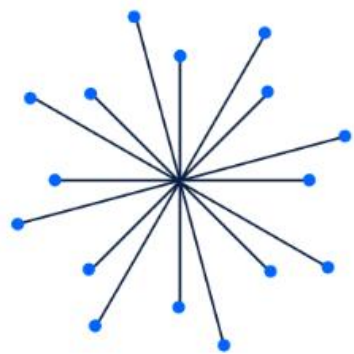
모바일시스템공학과
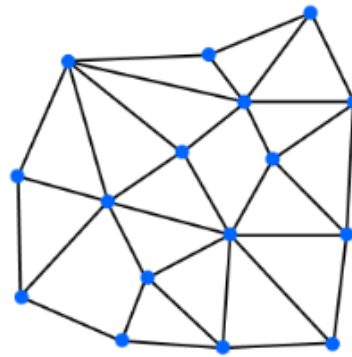이승재

# 0. How to find a paper

# 1. Introduction – Motivation

"As machine learning continues to gain popularity, sensitive user privacy has become an important issue, and FL has emerged, but it also has difficulties in ensuring complete privacy and accuracy."
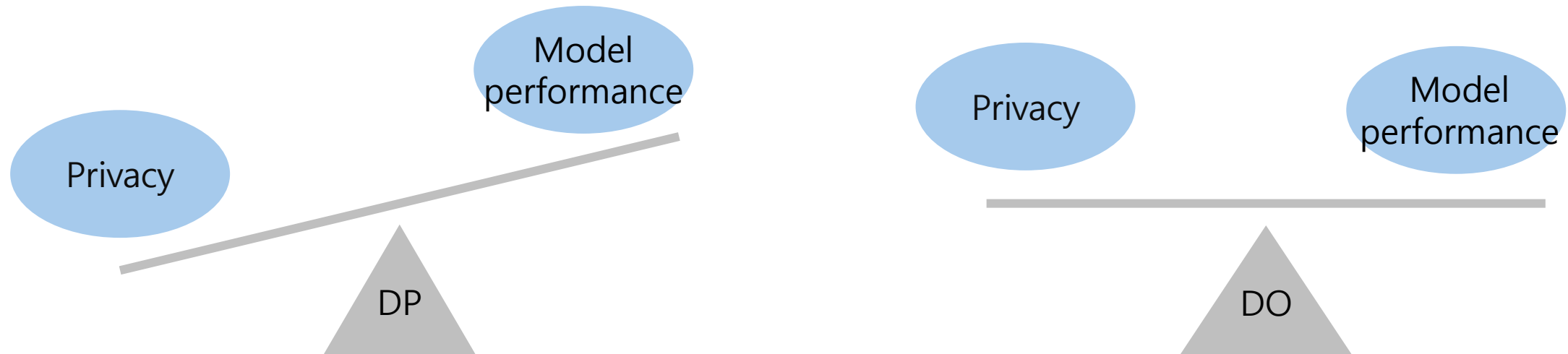


Centralized

Distributed
(FL)

+

Differential privacy
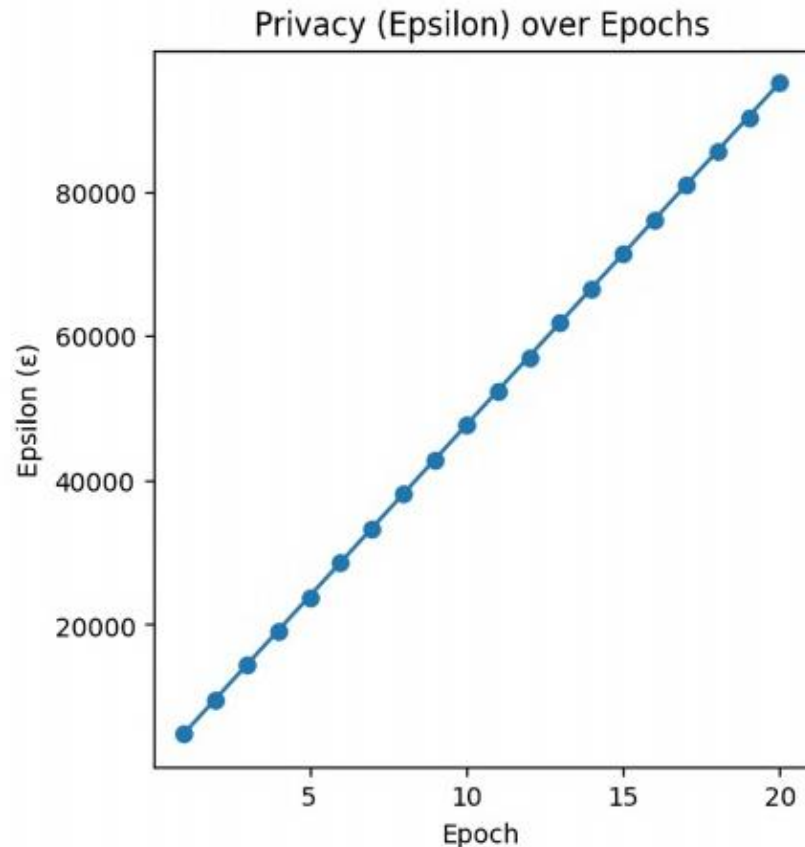(DP)

Data obfuscation
(DO)

!

# 1. Introduction – Motivation

"Existing methods such as DP can significantly degrade model performance by adding noise to data, making it difficult to balance privacy and accuracy, so new approaches for balanced models have been explored."

# 1. Introduction – Motivation

Privacy (Epsilon) over Epochs

Figure 7. Comparative analysis of Epsilon ($\epsilon$) vs. Epochs for FL-DP.

"DP not only has low model performance, but also reduces privacy protection strength as learning progresses."

| Training time ↑ | ε ↑ | Privacy guarantees ↓ |
| --- | --- | --- |

Data exposure risk when adding noise to data in DP

**Differential Privacy**

Data

Noise

$$\mathbb{P}[\mathcal{A}(D_1) \in S] \leq e^{\epsilon} \cdot \mathbb{P}[\mathcal{A}(D_2) \in S].$$

https://seewoo5.tistory.com/23

**Data Obfuscation**

Data

# 2. Background – More about DO

DO techniques can be used to secure sensitive information within the model because they make it difficult for attackers to interpret or understand the data, ensuring information is confidential. These methods have a few common techniques, including data masking, which involves replacing sensitive data with realistic but false information; encryption, which transforms data into a coded format requiring a key for decryption; and tokenization, where sensitive data elements are substituted with non-sensitive equivalents. Other methods include data shuffling, which rearranges entries in a database to hide connections. Perturbation adds noise or makes small changes to numerical data. Generalization reduces the detail of data, like changing specific ages to age ranges. Data swapping involves exchanging values between individual records. Additionally, nulling or deleting sensitive data replaces them with null values, making these techniques important for data protection.
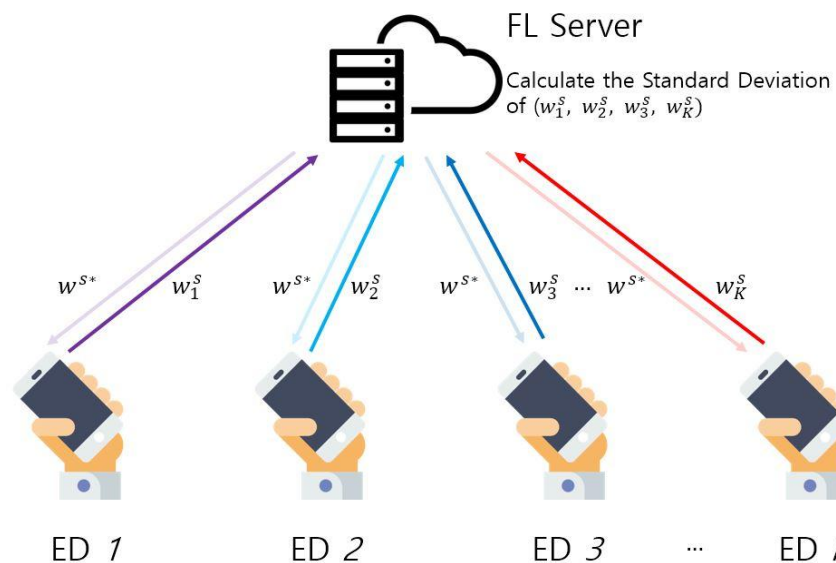
**Algorithm 1** FederatedAveraging. The $K$ clients are indexed by $k$; $B$ is the local minibatch size, $E$ is the number of local epochs, and $\eta$ is the learning rate.

**Server executes:**
  initialize $w_0$
  **for** each round $t = 1, 2, \ldots$ **do**
    $m \leftarrow \max(C \cdot K, 1)$
    $S_t \leftarrow$ (random set of $m$ clients)
    **for** each client $k \in S_t$ **in parallel do**
      $w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$
    $m_t \leftarrow \sum_{k \in S_t} n_k$
    $w_{t+1} \leftarrow \sum_{k \in S_t} \frac{n_k}{m_t} w_{t+1}^k$   *// Erratum[4]*

**ClientUpdate**$(k, w)$:   *// Run on client k*
  $\mathcal{B} \leftarrow$ (split $\mathcal{P}_k$ into batches of size $B$)
  **for** each local epoch $i$ from 1 to $E$ **do**
    **for** batch $b \in \mathcal{B}$ **do**
      $w \leftarrow w - \eta \nabla \ell(w; b)$
  return $w$ to server



"FL is a machine learning technique that enables multiple devices or systems to collaboratively train a model without sharing raw data. Instead of sending data to a central server, each device processes its data locally and shares only model updates (e.g., gradients) with the central server."
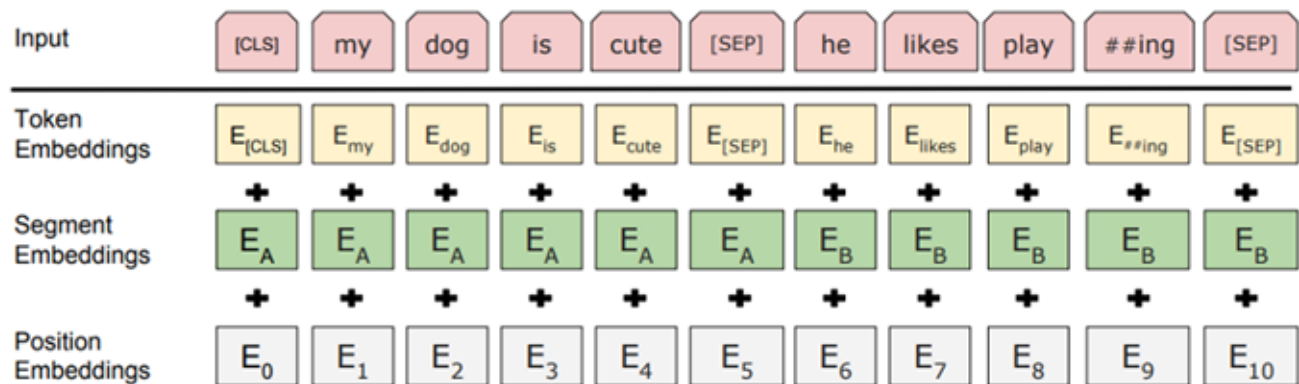


FL Server
Calculate the Standard Deviation of $(w_1^s, w_2^s, w_3^s, w_K^s)$

$w^{s*}$   $w_1^s$   $w^{s*}$   $w_2^s$   $w^{s*}$   $w_3^s \ldots w^{s*}$   $w_K^s$

ED *1*     ED *2*     ED *3*    ...    ED *K*

# 2. Background – BERT Configuration



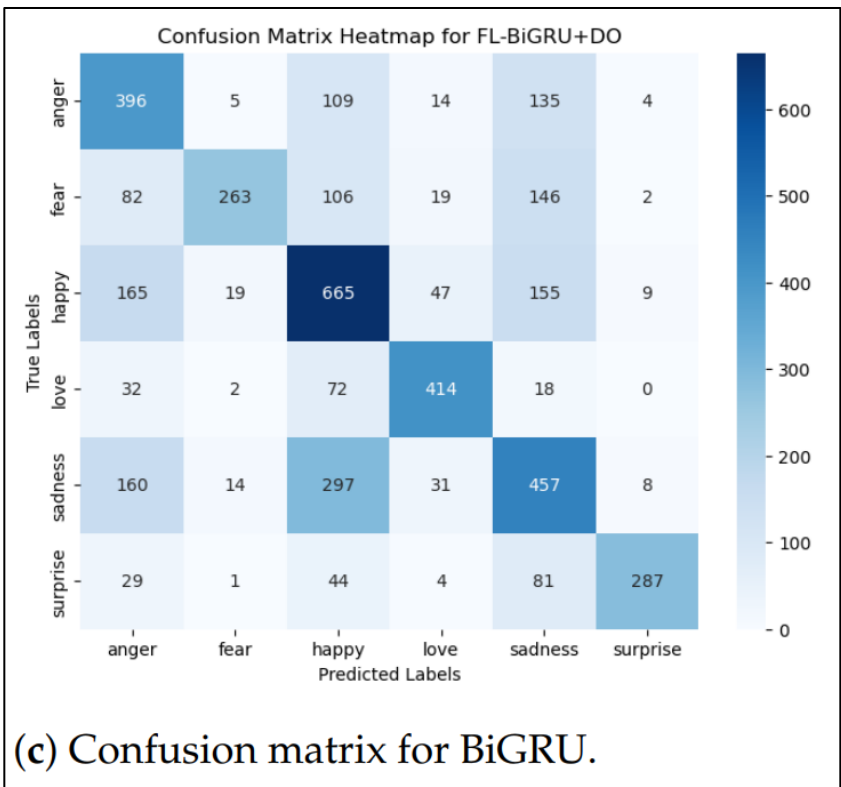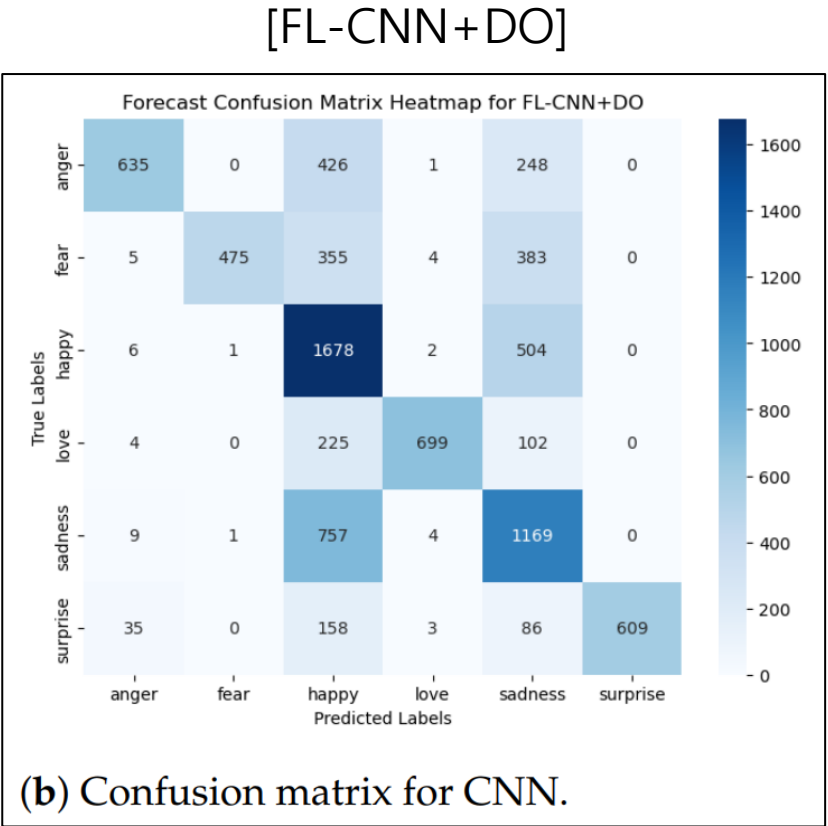| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

"Optimization factors for training the BERT model in this paper"

- Cross-Entropy Loss
- Learning Rate : $1 \times 10^{-5}$
- AdamW
- Weight Decay
- Learning Rate Scheduler
- Mixed-Precision Training
- Batch Size = 16

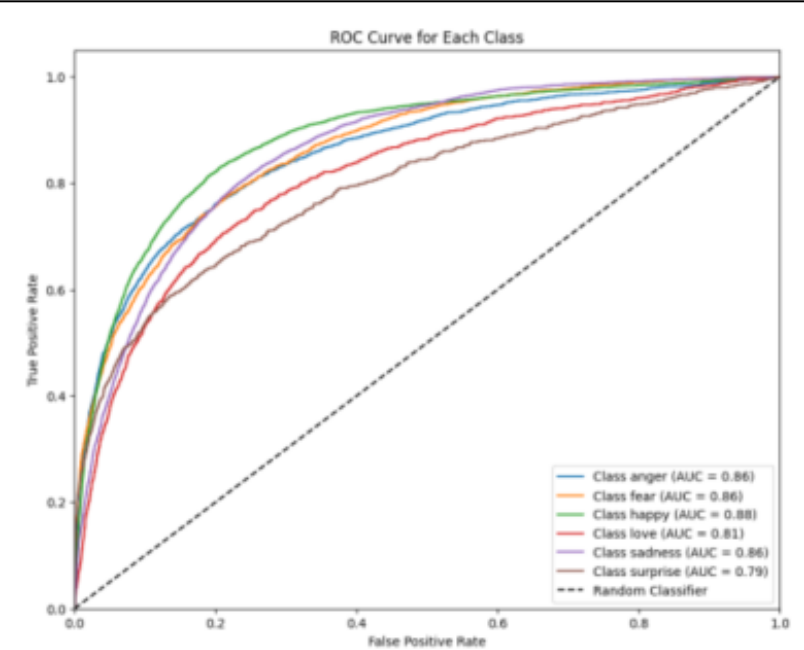# 3. Experimentation – Emotional Perspective



(a) Confusion matrix for BERT.

[FL-BERT+DO]

[FL-CNN+DO]



(b) Confusion matrix for CNN.



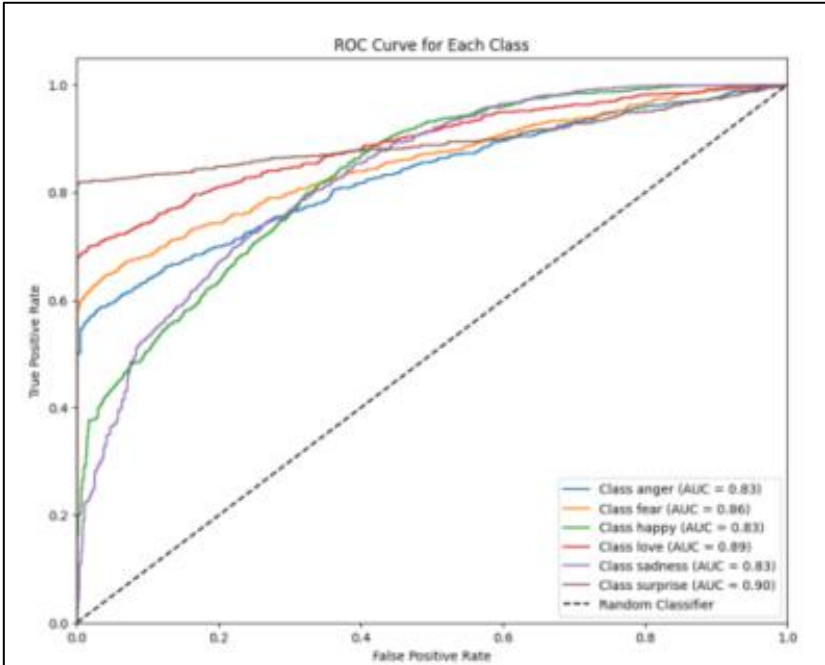(c) Confusion matrix for BiGRU.

[FL-BiGRU+DO]
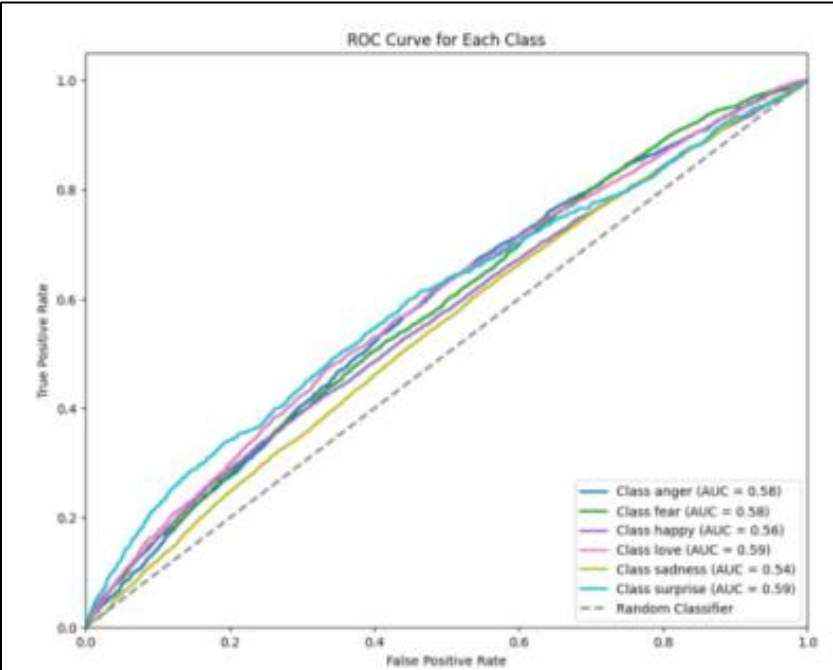
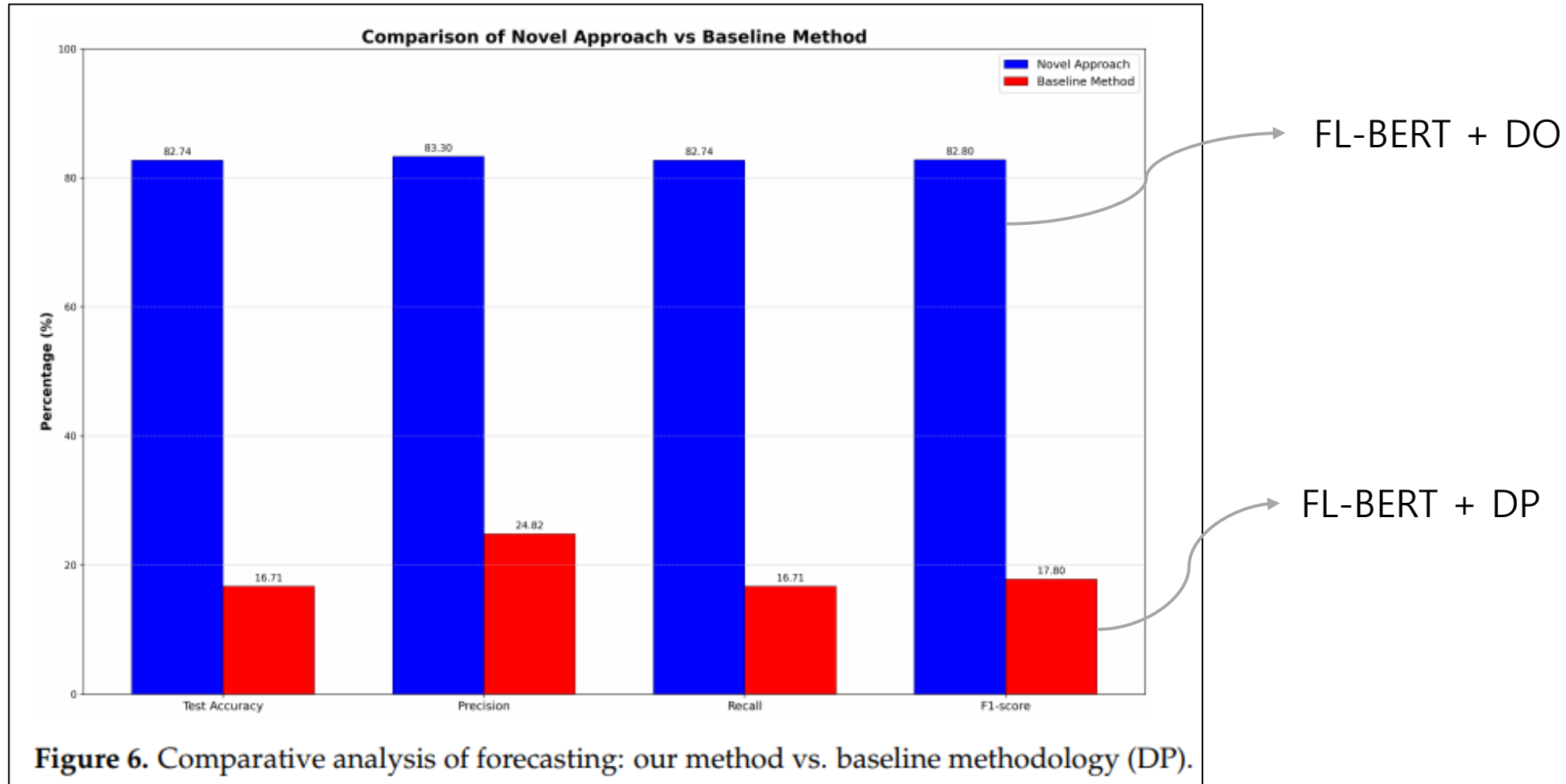# 3. Experimentation – Emotional Perspective



(a) ROC-AUC curve for BERT.

[FL-BERT+DO]

[FL-CNN+DO]

(b) ROC-AUC curve for CNN.

(c) ROC-AUC curve for BiGRU.

[FL-BiGRU+DO]

# 2. Experimentation – Emotional Perspective



**Figure 6.** Comparative analysis of forecasting: our method vs. baseline methodology (DP).

FL-BERT + DO

FL-BERT + DP

# 2. Experimentation - Privacy Perspective

**Table 4.** Privacy validation results for membership inference and linkage attacks.

| Attack Type | Model Type | AUC Score | Privacy Risk |
|---|---|---|---|
| **FL-BERT+DO** | | | |
| Membership Inference | Global BERT | 22.40% | Low |
| Membership Inference | Local BERT | 50.38% | Moderate |
| Linkage Attack | Individual Clients (Macro-Avg.) | 51.29% | Moderate |
| **FL-CNN+DO** | | | |
| Membership Inference | Global CNN | 37.36% | Low |
| Membership Inference | Local CNN | 50.95% | Moderate |
| Linkage Attack | Individual Clients (Macro-Avg.) | 50.72% | Moderate |
| **FL-BiGRU+DO** | | | |
| Membership Inference | Global BiGRU | 12.97% | Very Low |
| Membership Inference | Local BiGRU | 31.48% | Low |
| Linkage Attack | Individual Clients (Macro-Avg.) | 44.72% | Moderate |

It has the lowest risk to privacy, but it is not balanced with the model performance as shown in the previous experiment.

# 3. Future Directions

"Some prospects for future work with regard to data acquired for data obfuscation techniques include further exploration of the duality between privacy preservation and the usefulness of data. Future work might look into higher-level analogues of masking that would retain as much information as possible but would not compromise privacy.

Moreover, further research has to be directed to cognitive obscuring techniques by considering the characteristics of data and the context in which it will be used. Moreover, procedures for normalizing and subsequently verifying that masked data are still appropriate for the learning algorithms would be useful."