# Generative Federated Learning Approaches for Non-IID Data

Lohithsai Yadala Chanchu

**Supervisor:** Swier Garst    **Responsible Professor:** David Tax
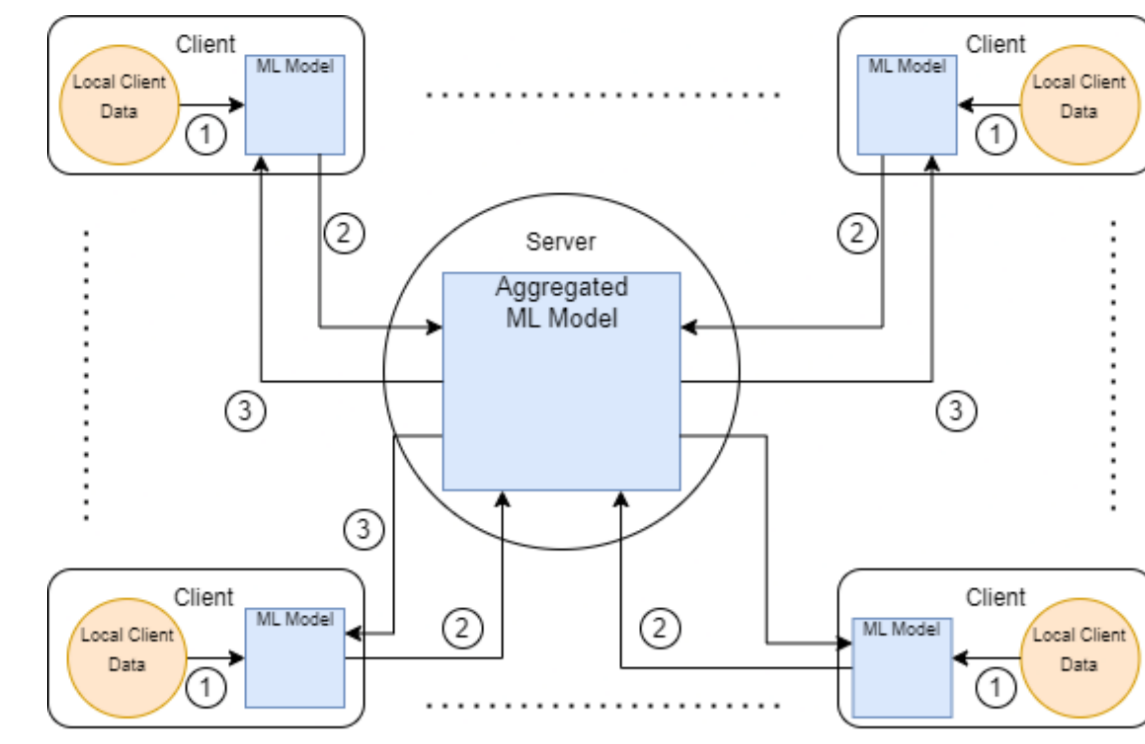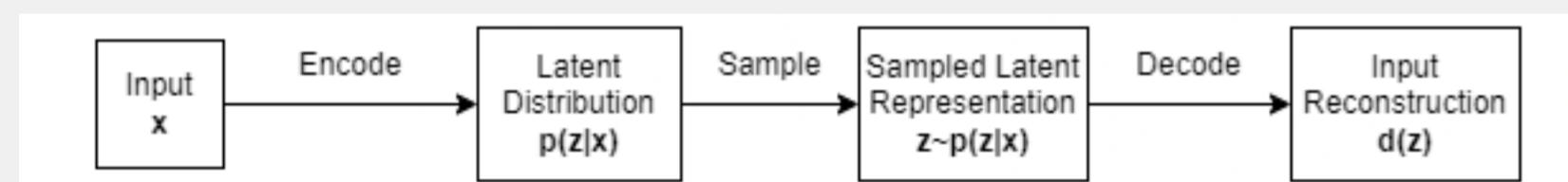
## 1. Introduction and Background



1. **What is Federated Learning (FL)?**: In FL systems, decentralized clients update a global model by training on local data. The central server coordinates the training process, collecting the locally trained models, and aggregating them to form an updated global model. This ensures local models can get the benefit of a bigger dataset while maintaining the privacy of their data.

2. **The Non-IID Problem**:
   Local client datasets not identically and independently distributed (IID) can significantly hamper global model performance due to the discrepancy between the global model and local model weights during training.

3. **Types of Solutions**:
   - **Data-based**: The data distribution is transformed in some way
   - **Framework based**: Focus on modifying the overall FL procedure. Some solution types include similarity clustering and knowledge distillation.

## 2. What are (C)VAEs



- VAEs produce a latent distribution instead of a point
- The loss function is a balance between the reconstruction term and the Kullback-Leibler distance term. This ensures that the sampled points from the latent distribution are encoded into "meaningful" points while ensuring continuity in the latent space.
- Conditional VAEs differ in that they produce data with pre-specified labels, whereas labels generated from vanilla VAEs are non-deterministic
- CVAEs ensure controlled label distribution, as vanilla VAEs trained in a federated setting do not produce the same distribution as the global dataset.
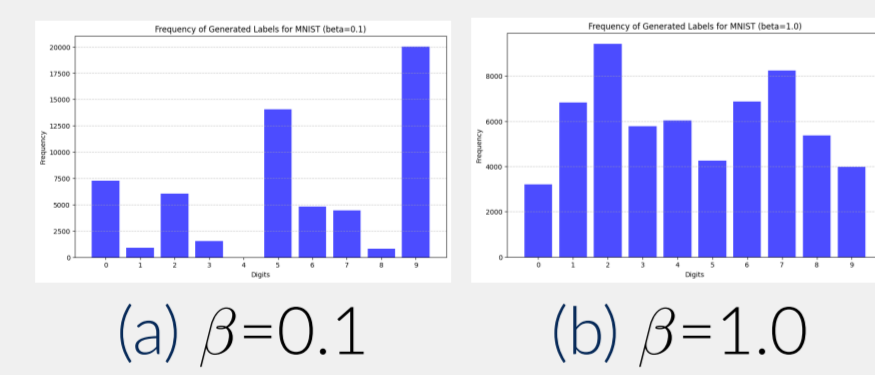


(a) $\beta$=0.1    (b) $\beta$=1.0

Figure 1. Influence of data imbalance on VAE-generated label distribution

## 3. Research Question

1. How much do generative techniques help federated performance as data imbalance increases?
2. To what extent does the choice of augmentation method affect federated performance?

### Hypothesis

- Augmentation based on local minority classes will lead to faster convergence and better accuracy metrics compared to no augmentation
- The less IID client datasets are, the greater the improvement across the tested frameworks and data-based approaches.

## 4. Methodology

**General**:

1. Dataset preparation
   - Partition the binarized MNIST dataset into different label skews using the Dirichlet distribution for 100 clients.
   - Vary the $\beta$ parameters: [0.1, 0.3, 0.5, 10.0, 100.0].
   - A smaller $\beta$ parameter results in greater imbalance, leading to increasingly skewed label distributions on clients.
2. CVAE is trained in a federated manner.
3. Trained CVAE is used to augment the dataset in various ways

**Uniform and Minority Imputation Methodology**:

1. Trained CVAE is used to augment $k$=[0, 100, 500, 1000] data points in the local dataset with the following imputation techniques:
   - **Uniform imputation:** Imputes all the labels uniformly, no matter what the initial local client label distribution is.
   - **Minority Imputation:** Takes the prior label distribution into account, imputing more synthetic data points with the minority labels; the bigger the difference between the most common label and the frequency of a certain label, the more that certain label is imputed.
2. The classifier is trained on the augmented local dataset in a federated manner.

**Finetuning Methodology**:

1. The CVAE generates a synthetic dataset which the central classifier uses to train itself locally on the server.
2. This pre-trained classifier is sent to each client.
3. The pre-trained classifier is further trained with original data from the local dataset.



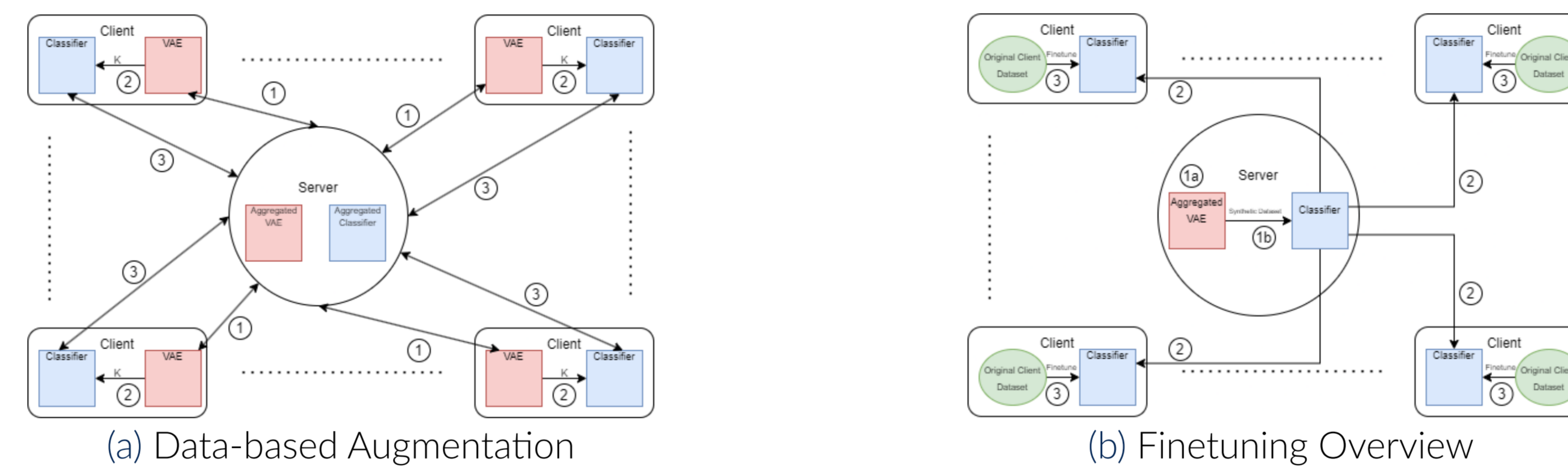(a) Data-based Augmentation    (b) Finetuning Overview

Figure 2. (a) Data-based Augmentation and (b) Finetuning Overview

## 6. Conclusion

- Combining original and synthetic data outperforms only original data in scenarios with imbalanced label distributions and only synthetic data on all $\beta$s.
- There is an optimal amount of synthetic data imputation for uniform imputation; beyond this, performance starts to degrade.
- Minority imputation seems to lead to better performance than uniform imputation across most levels of imbalances, with the most improvement for the most imbalanced data.

**Limitations**:

- **Absence of $\beta$-VAE:** Implementing a $\beta$-VAE could enhance the diversity of generative outputs.
- **Lack of Privacy Exploration:** Future research should investigate privacy-accuracy trade-offs with random noise in generated data points.
- **Focus on Limited Client Participation:** Future work should explore data augmentation in large federated learning systems for scalability and applicability.
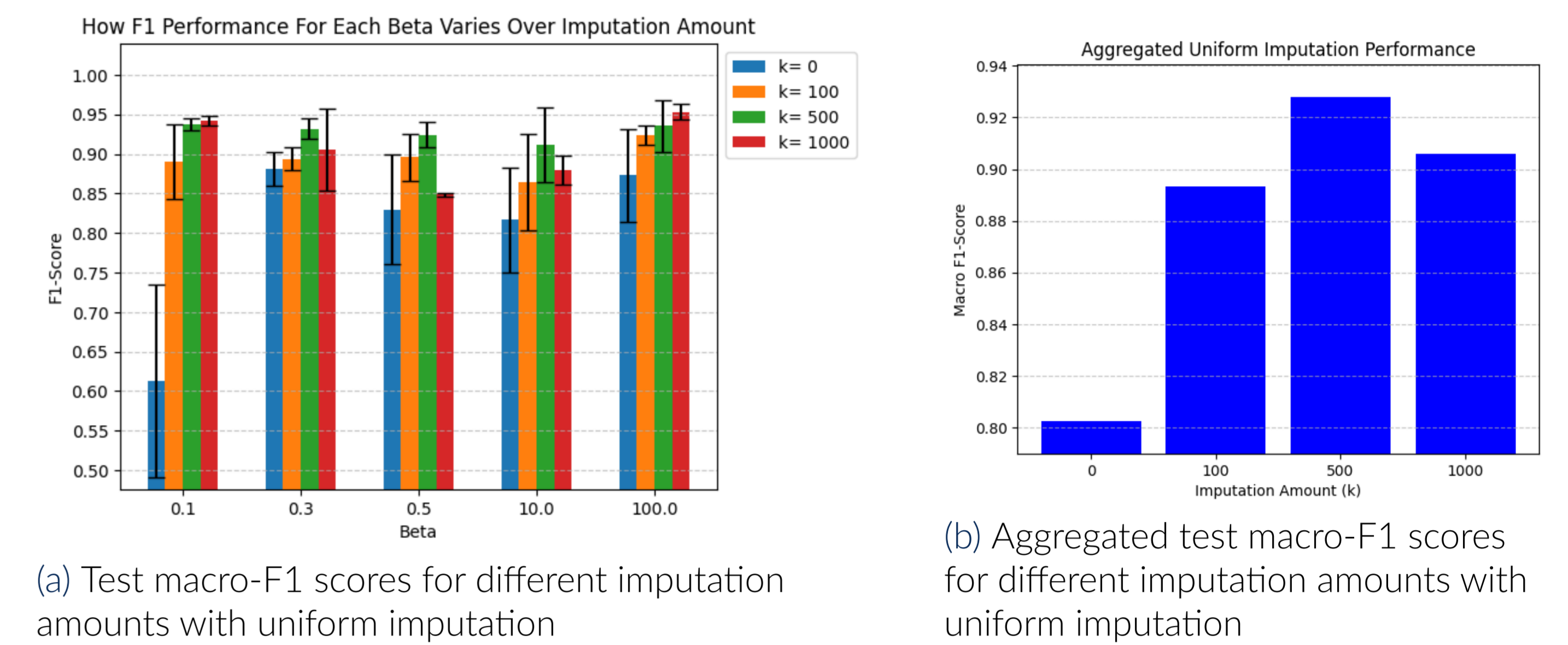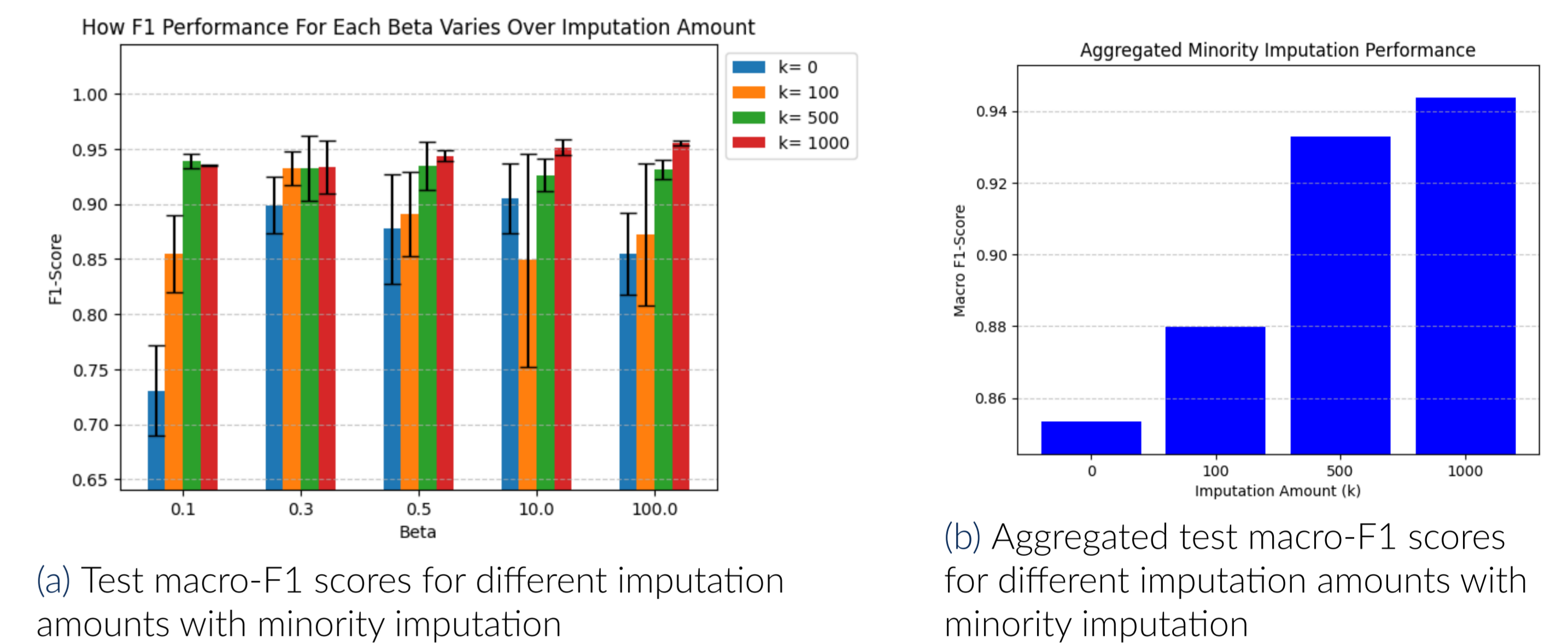
## 5. Results



(a) Test macro-F1 scores for different imputation amounts with uniform imputation

(b) Aggregated test macro-F1 scores for different imputation amounts with uniform imputation

Figure 3. Results for uniform imputation



(a) Test macro-F1 scores for different imputation amounts with minority imputation

(b) Aggregated test macro-F1 scores for different imputation amounts with minority imputation

Figure 4. Results for minority imputation



(a) Accuracies for classifier trained on synthetic MNIST data

(b) Accuracies for classifier trained on original MNIST data

(c) Test macro-F1 scores for pre-trained, finetuned, and vanilla federated strategies
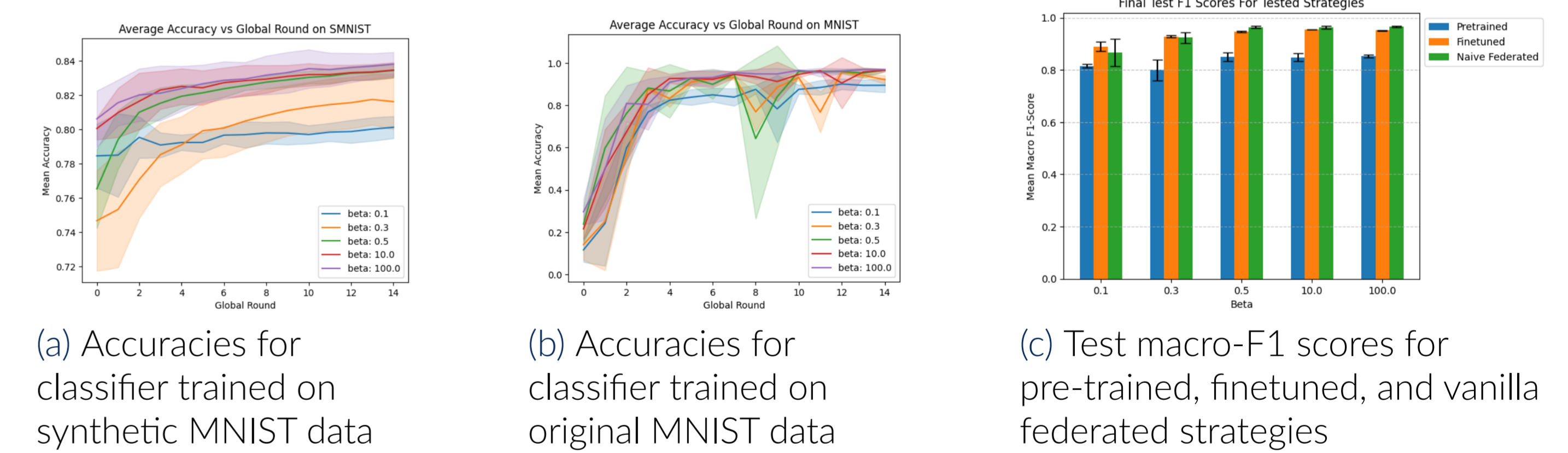
Figure 5. Results for imputation experiments

## References

[1] H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.
    Communication-efficient learning of deep networks from decentralized data.
    February 2016.

[2] Diletta Chiaro, Edoardo Prezioso, Michele Ianni, and Fabio Giampaolo.
    Fl-enhance: A federated learning framework for balancing non-iid data with augmented and shared compressed samples.
    *Information Fusion*, 98:101836, 2023.

[3] Taylor G. Smith and Jason M. White.
    Tgsmith61591/smrt: Handle class imbalance intelligently by using variational auto-encoders to generate synthetic observations of your minority class., 2017.

[4] Suman Ravuri and Oriol Vinyals.
    Classification accuracy score for conditional generative models.
    May 2019.