# Fall 2020 ECE M209AS Bake-off #2 Proposal

(Due: November 2, 11:59 PM)

## Instructions:

Main idea: **Human-AI Interaction**, broadly construed;

You will propose a self-defined project that is (1) enabled by AI, (2) interactive with human users, and (3) solving a domain-specific problem.

Please make a copy of this google document into your account and then fill the information below as well as answer the questions on the following pages (each question is on a new page). At the end, download as PDF and upload it here. Only one submission per team is required.

## Project Title:

GUI-GAN: Towards an interactive graphical framework for privacy-preserving artificial data synthesis and imputation using generative adversarial networks.

## Project Team:

Viacheslav Inderiakin; v.inderiakin.uk@gmail.com

Swapnil Sayan Saha; swapnilsayan@g.ucla.edu

**1. What are you trying to do in your project? What kind of AI is involved? How is it interactive with users? What domain-specific problems does it solve?**
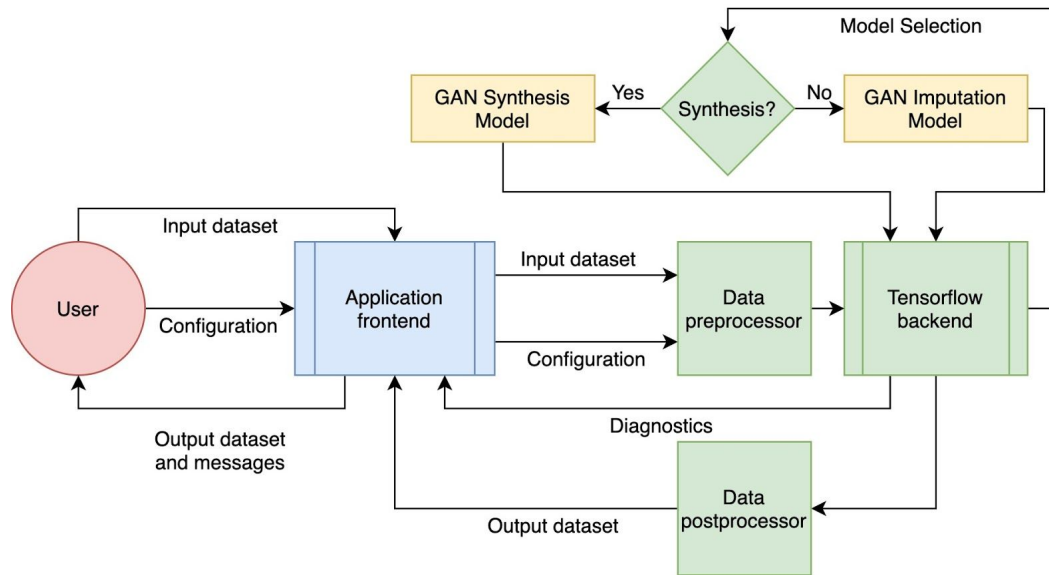
Your Response:

One of the major challenges in developing learning-enabled interactive systems is the absence of well-annotated large datasets, which can curtail the deployment performance of data-driven algorithms in the wild. Moreover, large-scale data analytics relying on mining of personal sensitive data lacks user adoption due to the prospects of privacy violation and potential of unearthing side-channel information from seemingly benign user data. As a result, there is a growing need for synthesizing artificial training datasets from limited training samples preserving required feature statistics for augmenting runtime effectiveness of machine-learning systems, while protecting user-privacy by eliminating side-channels. However, existing frameworks for artificial data synthesis are application-specific, uninteractive and require both computational and domain expertise, limiting their usage within a bounded-realm.

In this project, we design a real-time and interactive graphical user interface (GUI) framework for synthesizing large time-series datasets from moderately-sized input datasets using Generative Adversarial Networks (GANs) called GUI-GAN. Specifically, we illustrate a generalizable graphical pipeline powered by a deep-learning backend that allows non-experts to generate synthetic sensory data while parameterizing synthetic data statistics via graphical tools. In addition, the framework allows the user to symbolically specify rectification statistics of generated data, which is then fed to another generative deep-learning pipeline for data imputation based on user-corrections.

For evaluation, we show the application of our pipeline in the medical domain (where there is a lack of publicly-available datasets due to privacy redtapes) to generate privacy-preserving artificial human inertial motion and electrocardiogram (ECG) datasets from small samples. We believe that such a framework will pave the way for broad deployment of the power of artificial intelligence in eclectic domains (especially the medical domain) without requiring significant computational prowess, while certifying user confidence in privacy preservation for adopting data sharing.

**2. What is your proposed approach? Describe the planned system architecture, algorithms etc. Take as much space as needed, and include figures if necessary.**

Your Response:



The backend contains a data-preprocessor responsible for converting raw user data into tensor according to user requirements (e.g. window size, labels and statistical parameters). The framework then selects a synthesis GAN or an imputation GAN based on whether the user requires artificial dataset generation or generated dataset rectification (e.g. removing abnormal samples in temporal streams). The GANs are trained accordingly with input dataset using a Tensorflow backend, which yields either a synthetic dataset or corrected dataset based on user parameters. The backend is a black-box to the end-user, with the graphical application frontend bridging the gap between the "human" and the "machine". The user can interact with the framework through various graphical entities such as graphs, buttons and checkboxes.

To yield a highly generalizable time-series processor, we use domain-agnostic architectures for both GANs applicable to eclectic problem scenarios. For data synthesis, we evaluate an ensemble of two generative adversarial architectures:

- SenseGen - Generator: Stacked LSTM and Gaussian Mixture Model; Discriminator: LSTM
- TimeGAN - Generator: RNN; Discriminator: Bidirectional RNN and FFN, coupled with autoencoding components (embedding: RNN, recovery: FFN).

For data imputation, we feed the input data, mask metadata characterizing missing data and a random matrix to the generator network, while the discriminator uses the generator output (imputed data) coupled with a hint matrix to try and distinguish between observed and imputed components. This model is called GAIN.

**3. What if anything have you done thus far towards the project? This can include literature search (in which case provide a list of relevant websites, articles, papers etc. that you have identified).**

Your Response:

We have found several GAN frameworks that will be useful for the deep-learning backend of our proposed system:

1. Alzantot, Moustafa, Supriyo Chakraborty, and Mani Srivastava. "*Sensegen: A deep learning architecture for synthetic sensor data generation.*" 2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops). IEEE, 2017.

   Code: https://github.com/nesl/sensegen

   SenseGen is a generalizable GAN framework for generating privacy-preserving synthetic one-dimensional time-series datasets.

2. Yoon, Jinsung, James Jordon, and Mihaela van der Schaar. "*GAIN: Missing Data Imputation using Generative Adversarial Nets.*" International Conference on Machine Learning. 2018.

   Code: https://github.com/jsyoon0823/GAIN

   GAIN is a generative adversarial tool to impute missing data in temporal streams, significantly outperforming classical data imputation techniques.

3. Yoon, Jinsung, Daniel Jarrett, and Mihaela van der Schaar. "*Time-series generative adversarial networks.*" Advances in Neural Information Processing Systems. 2019

   Code: https://github.com/jsyoon0823/TimeGAN

   TimeGAN is a GAN framework for generating synthetic time-series datasets, with special stress put on following temporal dynamics of training datasets.

**4. What are the metrics of success/failure for your project?**

Your Response:

Several metrics of success include:

- Privacy preservation - misclassification rate of discriminator in synthesis network.
- RMSE of output dataset statistics with input dataset statistics or desired statistics.
- Discriminative and predictive scores for output dataset.
- RMSE and AUROC of data imputation network with varying missing data rates.
- Congeniality of data imputation network.
- (Optional) Classification accuracy using generated dataset in tertiary machine learning models.
- UI/UX scores:
    - SUS score (coupled with Shapiro Wilk, Friedman and Wilcoxon tests)
    - TLX score (coupled with Shapiro Wilk, Friedman and Wilcoxon tests)
    - General feedback and suggestions

The metrics will be extracted for two cases in the medical domain:

- ECG dataset
- Human motion dataset from inertial sensors (e.g. Accelerometer)

**5. What items do you need or plan to use beyond Hexiwear, RPiZW, Smartphone, and Laptop? Please note that irrespective of whether the item is purchased by me, by you, or is already here with me or you, the total pre-tax/pre-shipping cost of all items must not exceed $25 * # of students in the team. It is okay if you don't yet know everything that you would need, but describe to the best of your ability.**

Your Response: We do not need anything beyond our laptops.

| | |
|---|---|
| Item Description | |
| Manufacturer | |
| Manufacturer Part # | |
| URL to item on manufacturer website | |
| Vendor | |
| Vendor Part # | |
| URL to item on vendor website | |
| Price | |
| Quantity | |

[Please replicate the above table as needed for every item your project needs].

Total Cost = <please fill>