

Argoverse'd :

Motion Forecasting with Deep Learning Trained on a New Dataset

Hardy Bright, Robyn Logan, Shane Parr, Ye Wang

Devpost link: <https://devpost.com/software/vadimaniacs-heresy#updates>

Github link: <https://github.com/Robyn42/Vadimaniac-Code>

Introduction

Motion forecasting is a multivariate time-series regression task in which the goal is to predict the future trajectories of moving objects. By analyzing the trajectory of objects in labeled data, algorithms can be trained to accurately predict the future trajectory of new objects based on their current positions and headings. This allows for autonomous vehicles to be proactive in their object avoidance instead of older reactive technologies used in driver-assist technologies.

We are training and assessing the robustness of deep learning models on the Argoverse 2 motion prediction dataset^[1], which consists of vehicles driving on public roads in six US cities (Austin, Detroit, Miami, Pittsburgh, Palo Alto, and Washington D.C.) to provide a robust set of data from across the country. For the purpose of motion forecasting, the data was processed into 2-dimensional maps to simplify the training and implementation of motion forecasting algorithms. The data was also labeled to allow the algorithms to differentiate objects as various vehicles and beings to further improve forecasting accuracy.

Our team has developed multiple model architectures to date including stacked long short-term memory (LSTM), multi-module gated recurrent units (GRU) in an observe/predict configuration, single module GRU, and multi-layer perceptron (MLP). An MLP was selected to serve as a good baseline with parallel structure that might align with the timeseries data, the GRUs were selected because of their improved speed and efficiency when compared to LSTMs, and the stacked LSTM was selected as a counterpoint to the “stacked” observation/prediction model. The models trained in this paper have been shown to perform better than models on Argoverse 1's public leaderboards in average displacement error (ADE) metrics and comparably in others.

Methodology

The Dataset

Argoverse 2 is a multimodal, multivariate dataset produced by Argo AI^[2]. It builds upon the foundation of Argoverse 1 by having four major components: sensor, lidar, motion forecasting and map change datasets were compiled to be the state-of-the-art for autonomous vehicle development. Our project makes use of the motion forecasting part of the dataset, which is portioned into train, validation and test datasets together totaling over 200 Gigabytes (GB). Over the course of development for our various models, we used each of the three for cross-validation. Our final training and evaluation was done with the train and validation set respectfully.

```
observed: bool
track_id: string
object_type: string
object_category: int64
timestep: int64
position_x: double
position_y: double
heading: double
velocity_x: double
velocity_y: double
scenario_id: string
start_timestamp: double
end_timestamp: double
num_timestamps: int64
focal_track_id: string
city: string
=====
observed: [[true,true,true,true,true,true,true,true,true,true,...,false,false,false,false,false,false,false,false,false]]
track_id: ["77543","77543","77543","77543","77543","77543","77543","77543","77543","77543","77543",...,"AV","AV","AV","AV","AV","AV","AV","AV","AV","AV"]
object_type: ["vehicle","vehicle","vehicle","vehicle","vehicle","vehicle","vehicle","vehicle","vehicle","vehicle",...,"vehicle","vehicle","vehicle"]
object_category: [[1,1,1,1,1,1,1,1,1,1,...,1,1,1,1,1,1,1,1,1,1]]
timestep: [[0,1,2,3,4,5,6,7,8,9,...,100,101,102,103,104,105,106,107,108,109]]
position_x: [[-11.823673649007935,-12.39320578754739,-13.080372362384209,-13.880497176060121,-14.773626809536655,-15.74815042844231,-16.7757018228538
position_y: [[-567.4023974854001,-567.2253156611156,-567.0096938093519,-566.7544103059896,-566.4653864371573,-566.1455193078166,-565.8032455480408,-565.5111111111111,-565.2222222222222,-564.9333333333333,-564.6444444444444,-564.3555555555555,-564.0666666666667,-563.7777777777778,-563.4888888888889,-563.1999999999999,-562.9111111111111,-562.6222222222222,-562.3333333333333,-562.0444444444444,-561.7555555555556,-561.4666666666667,-561.1777777777778,-560.8888888888889,-560.5999999999999,-560.3111111111111,-560.0222222222222,-559.7333333333333,-559.4444444444444,-559.1555555555556,-558.8666666666667,-558.5777777777778,-558.2888888888889,-557.9999999999999,-557.7111111111111,-557.4222222222222,-557.1333333333333,-556.8444444444444,-556.5555555555556,-556.2666666666667,-555.9777777777778,-555.6888888888889,-555.3999999999999,-555.1111111111111,-554.8222222222222,-554.5333333333333,-554.2444444444444,-553.9555555555556,-553.6666666666667,-553.3777777777778,-553.0888888888889,-552.7999999999999,-552.5111111111111,-552.2222222222222,-551.9333333333333,-551.6444444444444,-551.3555555555556,-551.0666666666667,-550.7777777777778,-550.4888888888889,-550.1999999999999,-549.9111111111111,-549.6222222222222,-549.3333333333333,-549.0444444444444,-548.7555555555556,-548.4666666666667,-548.1777777777778,-547.8888888888889,-547.5999999999999,-547.3111111111111,-547.0222222222222,-546.7333333333333,-546.4444444444444,-546.1555555555556,-545.8666666666667,-545.5777777777778,-545.2888888888889,-544.9999999999999,-544.7111111111111,-544.4222222222222,-544.1333333333333,-543.8444444444444,-543.5555555555556,-543.2666666666667,-542.9777777777778,-542.6888888888889,-542.3999999999999,-542.1111111111111,-541.8222222222222,-541.5333333333333,-541.2444444444444,-540.9555555555556,-540.6666666666667,-540.3777777777778,-540.0888888888889,-539.7999999999999,-539.5111111111111,-539.2222222222222,-538.9333333333333,-538.6444444444444,-538.3555555555556,-538.0666666666667,-537.7777777777778,-537.4888888888889,-537.1999999999999,-536.9111111111111,-536.6222222222222,-536.3333333333333,-536.0444444444444,-535.7555555555556,-535.4666666666667,-535.1777777777778,-534.8888888888889,-534.5999999999999,-534.3111111111111,-534.0222222222222,-533.7333333333333,-533.4444444444444,-533.1555555555556,-532.8666666666667,-532.5777777777778,-532.2888888888889,-531.9999999999999,-531.7111111111111,-531.4222222222222,-531.1333333333333,-530.8444444444444,-530.5555555555556,-530.2666666666667,-530.0777777777778,-529.7888888888889,-529.4999999999999,-529.2111111111111,-528.9222222222222,-528.6333333333333,-528.3444444444444,-528.0555555555556,-527.7666666666667,-527.4777777777778,-527.1888888888889,-526.8999999999999,-526.6111111111111,-526.3222222222222,-526.0333333333333,-525.7444444444444,-525.4555555555556,-525.1666666666667,-524.8777777777778,-524.5888888888889,-524.2999999999999,-524.0111111111111,-523.7222222222222,-523.4333333333333,-523.1444444444444,-522.8555555555556,-522.5666666666667,-522.2777777777778,-521.9888888888889,-521.6999999999999,-521.4111111111111,-521.1222222222222,-520.8333333333333,-520.5444444444444,-520.2555555555556,-520.0666666666667,-519.7777777777778,-519.4888888888889,-519.1999999999999,-518.9111111111111,-518.6222222222222,-518.3333333333333,-518.0444444444444,-517.7555555555556,-517.4666666666667,-517.1777777777778,-516.8888888888889,-516.5999999999999,-516.3111111111111,-516.0222222222222,-515.7333333333333,-515.4444444444444,-515.1555555555556,-514.8666666666667,-514.5777777777778,-514.2888888888889,-513.9999999999999,-513.7111111111111,-513.4222222222222,-513.1333333333333,-512.8444444444444,-512.5555555555556,-512.2666666666667,-511.9777777777778,-511.6888888888889,-511.3999999999999,-511.1111111111111,-510.8222222222222,-510.5333333333333,-510.2444444444444,-510.0555555555556,-509.7666666666667,-509.4777777777778,-509.1888888888889,-508.8999999999999,-508.6111111111111,-508.3222222222222,-508.0333333333333,-507.7444444444444,-507.4555555555556,-507.1666666666667,-506.8777777777778,-506.5888888888889,-506.2999999999999,-506.0111111111111,-505.7222222222222,-505.4333333333333,-505.1444444444444,-504.8555555555556,-504.5666666666667,-504.2777777777778,-503.9888888888889,-503.6999999999999,-503.4111111111111,-503.1222222222222,-502.8333333333333,-502.5444444444444,-502.2555555555556,-501.9666666666667,-501.6777777777778,-501.3888888888889,-501.0999999999999,-500.8111111111111,-500.5222222222222,-500.2333333333333,-500.0444444444444,-499.7555555555556,-499.4666666666667,-499.1777777777778,-498.8888888888889,-498.5999999999999,-498.3111111111111,-498.0222222222222,-497.7333333333333,-497.4444444444444,-497.1555555555556,-496.8666666666667,-496.5777777777778,-496.2888888888889,-495.9999999999999,-495.7111111111111,-495.4222222222222,-495.1333333333333,-494.8444444444444,-494.5555555555556,-494.2666666666667,-493.9777777777778,-493.6888888888889,-493.3999999999999,-493.1111111111111,-492.8222222222222,-492.5333333333333,-492.2444444444444,-491.9555555555556,-491.6666666666667,-491.3777777777778,-491.0888888888889,-490.7999999999999,-490.5111111111111,-490.2222222222222,-489.9333333333333,-489.6444444444444,-489.3555555555556,-489.0666666666667,-488.7777777777778,-488.4888888888889,-488.1999999999999,-487.9111111111111,-487.6222222222222,-487.3333333333333,-487.0444444444444,-486.7555555555556,-486.4666666666667,-486.1777777777778,-485.8888888888889,-485.5999999999999,-485.3111111111111,-485.0222222222222,-484.7333333333333,-484.4444444444444,-484.1555555555556,-483.8666666666667,-483.5777777777778,-483.2888888888889,-482.9999999999999,-482.7111111111111,-482.4222222222222,-482.1333333333333,-481.8444444444444,-481.5555555555556,-481.2666666666667,-480.9777777777778,-480.6888888888889,-480.3999999999999,-480.1111111111111,-479.8222222222222,-479.5333333333333,-479.2444444444444,-478.9555555555556,-478.6666666666667,-478.3777777777778,-478.0888888888889,-477.7999999999999,-477.5111111111111,-477.2222222222222,-476.9333333333333,-476.6444444444444,-476.3555555555556,-476.0666666666667,-475.7777777777778,-475.4888888888889,-475.1999999999999,-474.9111111111111,-474.6222222222222,-474.3333333333333,-474.0444444444444,-473.7555555555556,-473.4666666666667,-473.1777777777778,-472.8888888888889,-472.5999999999999,-472.3111111111111,-472.0222222222222,-471.7333333333333,-471.4444444444444,-471.1555555555556,-470.8666666666667,-470.5777777777778,-470.2888888888889,-469.9999999999999,-469.7111111111111,-469.4222222222222,-469.1333333333333,-468.8444444444444,-468.5555555555556,-468.2666666666667,-467.9777777777778,-467.6888888888889,-467.3999999999999,-467.1111111111111,-466.8222222222222,-466.5333333333333,-466.2444444444444,-465.9555555555556,-465.6666666666667,-465.3777777777778,-465.0888888888889,-464.7999999999999,-464.5111111111111,-464.2222222222222,-463.9333333333333,-463.6444444444444,-463.3555555555556,-463.0666666666667,-462.7777777777778,-462.4888888888889,-462.1999999999999,-461.9111111111111,-461.6222222222222,-461.3333333333333,-461.0444444444444,-460.7555555555556,-460.4666666666667,-460.1777777777778,-459.8888888888889,-459.5999999999999,-459.3111111111111,-459.0222222222222,-458.7333333333333,-458.4444444444444,-458.1555555555556,-457.8666666666667,-457.5777777777778,-457.2888888888889,-456.9999999999999,-456.7111111111111,-456.4222222222222,-456.1333333333333,-455.8444444444444,-455.5555555555556,-455.2666666666667,-454.9777777777778,-454.6888888888889,-454.3999999999999,-454.1111111111111,-453.8222222222222,-453.5333333333333,-453.2444444444444,-452.9555555555556,-452.6666666666667,-452.3777777777778,-452.0888888888889,-451.7999999999999,-451.5111111111111,-451.2222222222222,-450.9333333333333,-450.6444444444444,-450.3555555555556,-450.0666666666667,-449.7777777777778,-449.4888888888889,-449.1999999999999,-448.9111111111111,-448.6222222222222,-448.3333333333333,-448.0444444444444,-447.7555555555556,-447.4666666666667,-447.1777777777778,-446.8888888888889,-446.5999999999999,-446.3111111111111,-446.0222222222222,-445.7333333333333,-445.4444444444444,-445.1555555555556,-444.8666666666667,-444.5777777777778,-444.2888888888889,-443.9999999999999,-443.7111111111111,-443.4222222222222,-443.1333333333333,-442.8444444444444,-442.5555555555556,-442.2666666666667,-441.9777777777778,-441.6888888888889,-441.3999999999999,-441.1111111111111,-440.8222222222222,-440.5333333333333,-440.2444444444444,-439.9555555555556,-439.6666666666667,-439.3777777777778,-439.0888888888889,-438.7999999999999,-438.5111111111111,-438.2222222222222,-437.9333333333333,-437.6444444444444,-437.3555555555556,-437.0666666666667,-436.7777777777778,-436.4888888888889,-436.1999999999999,-435.9111111111111,-435.6222222222222,-435.3333333333333,-435.0444444444444,-434.7555555555556,-434.4666666666667,-434.1777777777778,-433.8888888888889,-433.5999999999999,-433.3111111111111,-433.0222222222222,-432.7333333333333,-432.4444444444444,-432.1555555555556,-431.8666666666667,-431.5777777777778,-431.2888888888889,-430.9999999999999,-430.7111111111111,-430.4222222222222,-430.1333333333333,-429.8444444444444,-429.5555555555556,-429.2666666666667,-428.9777777777778,-428.6888888888889,-428.3999999999999,-428.1111111111111,-427.8222222222222,-427.5333333333333,-427.2444444444444,-426.9555555555556,-426.6666666666667,-426.3777777777778,-426.0888888888889,-425.7999999999999,-425.5111111111111,-425.2222222222222,-424.9333333333333,-424.6444444444444,-424.3555555555556,-424.0666666666667,-423.7777777777778,-423.4888888888889,-423.1999999999999,-422.9111111111111,-422.6222222222222,-422.3333333333333,-422.0444444444444,-421.7555555555556,-421.4666666666667,-421.1777777777778,-420.8888888888889,-420.5999999999999,-420.3111111111111,-420.0222222222222,-419.7333333333333,-419.4444444444444,-419.1555555555556,-418.8666666666667,-418.5777777777778,-418.2888888888889,-417.9999999999999,-417.7111111111111,-417.4222222222222,-417.1333333333333,-416.8444444444444,-416.5555555555556,-416.2666666666667,-415.9777777777778,-415.6888888888889,-415.3999999999999,-415.1111111111111,-414.8222222222222,-414.5333333333333,-414.2444444444444,-413.9555555555556,-413.6666666666667,-413.3777777777778,-413.0888888888889,-412.7999999999999,-412.5111111111111,-412.2222222222222,-411.9333333333333,-411.6444444444444,-411.3555555555556,-411.0666666666667,-410.7777777777778,-410.4888888888889,-410.1999999999999,-409.9111111111111,-409.6222222222222,-409.3333333333333,-409.0444444444444,-408.7555555555556,-408.4666666666667,-408.1777777777778,-407.8888888888889,-407.5999999999999,-407.3111111111111,-407.0222222222222,-406.7333333333333,-406.4444444444444,-406.1555555555556,-405.8666666666667,-405.5777777777778,-405.2888888888889,-404.9999999999999,-404.7111111111111,-404.4222222222222,-404.1333333333333,-403.8444444444444,-403.5555555555556,-403.2666666666667,-402.9777777777778,-402.6888888888889,-402.3999999999999,-402.1111111111111,-401.8222222222222,-401.5333333333333,-401.2444444444444,-400.9555555555556,-400.6666666666667,-400.3777777777778,-400.0888888888889,-399.7999999999999,-399.5111111111111,-399.2222222222222,-398.9333333333333,-398.6444444444444,-398.3555555555556,-398.0666666666667,-397.7777777777778,-397.4888888888889,-397.1999999999999,-396.9111111111111,-396.6222222222222,-396.3333333333333,-396.0444444444444,-395.7555555555556,-395.4666666666667,-395.1777777777778,-394.8888888888889,-394.5999999999999,-394.3111111111111,-394.0222222222222,-393.7333333333333,-393.4444444444444,-393.1555555555556,-392.8666666666667,-392.5777777777778,-392.2888888888889,-391.9999999999999,-391.7111111111111,-391.4222222222222,-391.1333333333333,-390.8444444444444,-390.5555555555556,-390.2666666666667,-389.9777777777778,-389.6888888888889,-389.3999999999999,-389.1111111111111,-388.8222222222222,-388.5333333333333,-388.2444444444444,-387.9555555555556,-387.6666666666667,-387.3777777777778,-387.0888888888889,-386.7999999999999,-386.5111111111111,-386.2222222222222,-385.9333333333333,-385.6444444444444,-385.3555555555556,-385.0666666666667,-384.7777777777778,-384.4888888888889,-384.1999999999999,-383.9111111111111,-383.6222222222222,-383.3333333333333,-383.0444444444444,-382.7555555555556,-382.4666666666667,-382.1777777777778,-381.8888888888889,-381.5999
```

Data preprocessing

Beginning our regression analysis, we chose to focus our attention on the data's continuous variables. Specifically, the `position_x`, `position_y`, `heading`, `velocity_x`, and `velocity_y` information. As stated earlier, we did not detect any missing data so imputation of values was not necessary. Our early models made use of the unaltered data while later analysis included normalization of the feature values to values between 0 and 1. As noted earlier, each scenario contains data on multiple vehicles. Although the main agents and many of the surrounding ones are recorded for the full duration, the dataset contains a number of vehicle observations that are not a full 11 seconds long. During initial model exploration, we decided to focus on data that represented vehicle behavior for a full 110 timesteps. Any time series less than that target number was disregarded; the disregarded datapoints represent 1% of the dataset. The decision was motivated by considering how analogous this numerical time series data is to natural language prediction tasks. In essence, we considered each full 110 timestep series of observations to be a sentence. By equating the observations to words in a sentence, the features then become the letters composing the word. We thus considered employing model architectures that excel at sequence to sequence translation, such as LSTM and GRU. Models such as these generally take windows of data as input into the modules. Adhering to only full length, 110 timestep sequences allowed us to treat the dataset as containing "sentences" of equal length that are easily divisible by window sizes of 10, where each window is composed of 10 "words".

To augment the data discussed earlier, three *social features* were computed using a helper function. That function uses the available data to compute the minimum distance (`MIN_DISTANCE_FRONT`), minimum distance back (`MIN_DISTANCE_BACK`) and the number of neighbors (`NUM_NEIGHBORS`) for the main agent of each scenario. These new features are concatenated along with the previous 16 feature variables and all of the timesteps are specified for the main agent vehicle only. The calculation of these new features also afforded some file-size advantages. The process creates new dataset files for train, validation and test sets that are significantly smaller and in Comma Separated Value (CSV) format. For example the training dataset is roughly 80 GB before processing and 1.5 GB afterwards. Once the new files are created, the original dataset can be disregarded.

Prediction Task

The task was to predict the trajectory of the main agent vehicle given the observed x and y positional data combined with the computed "social features" outlined above.

Model architecture

Over the course of the project, our team approached the prediction task on this new dataset through the development of 4 different model architectures. Our Observe 80/predict 30^[3], Stacked LSTM^[3, 4], single module GRU, and multi-layer perceptron models^[5] constitute different

conceptualizations of data flow to effectively address the multi-modal nature of our aggregated information.

All of the models we developed during the course of this project share some common hyperparameters and features:

Batch Size	128
Learning Rate	1e-3
Leaky ReLU Alpha	3e-1
Dropout Rate	3e-2
Hidden Layer Dimension	100
GRU or LSTM units	128

Figure 2. Shared baseline hyperparameter values.

Dropout was used during the training phase between each of the hidden dense layers. Each model utilizes Mean Squared Error (MSE) for its loss function. Weight optimization is achieved through the Adam optimizer.

Multi-layer perceptron model (MLP)

Our initial investigation began with this architecture where we have an input layer, one hidden layer and an output layer. The first two layers use Leaky ReLU as their activation function. During the training phase of the model, it uses Dropout between the layers.

It should be noted that this model was trained with a small portion of the dataset to address early storage capacity and computational concerns. The observations made here may have been influenced by that but what we experienced served as a foundation for the project's later work.

A significant portion of this stage of development involved cultivating an understanding of the multivariate nature of the data and how it should be prepared for prediction accuracy. The n-gram concept gave us an approach that worked well. Again, relating time series data to words and sentences, we created "5-grams". Each observation is a timestep and therefore a "word". Four "words" would go through the model to predict the fifth. We'd like to reiterate here that each "word" is composed of the feature variables (position_x, position_y, heading, velocity_x, velocity_y, MIN_DISTANCE_FRONT, MIN_DISTANCE_BACK, and NUM_NEIGHBORS). The model was tasked with predicting all eight of the variables given as input.

This model was taken a step further when we altered the makeup of the true values. Here we calculated the difference between the values of the fourth and fifth timesteps. Using these difference measures improved both training and testing loss significantly.

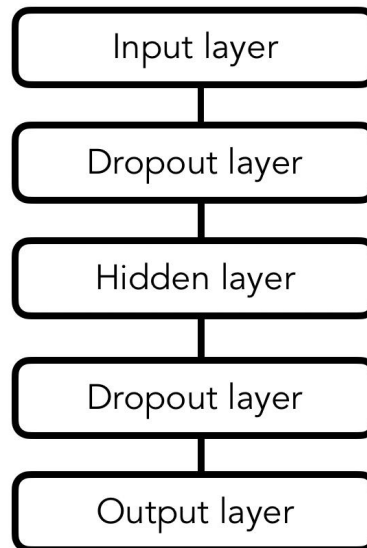


Figure 3. multi-layer perceptron model architectural diagram.

Single GRU module model

This model has a single GRU unit which processes all of its input and outputs it into a series of 2 dense layers. As with the previous model discussed, it utilizes LeakyReLU for layer activations along with Dropout during the training phase.

It was developed under the same limited data context as the multi-layer perceptron therefore the observations we gained from it may not be accurate of its full potential. It is included here because of its notable advantages and disadvantages to the multi-layer perceptron. We observed that the multi-layer perceptron reached a low level of loss with less epochs than this model. Conversely, the GRU's prediction accuracy appeared to be better than the MLP when given more epochs to train.

An explanation of this may reside in the way data is structured within n-grams. Because of its construction, the n-gram produces more data per scenario file in comparison to the windows of the GRU data preparation. That leads to more data being available during each epoch of training for the MLP. The GRU in effect needs to "catch up" with more epochs due to less available data per epoch.

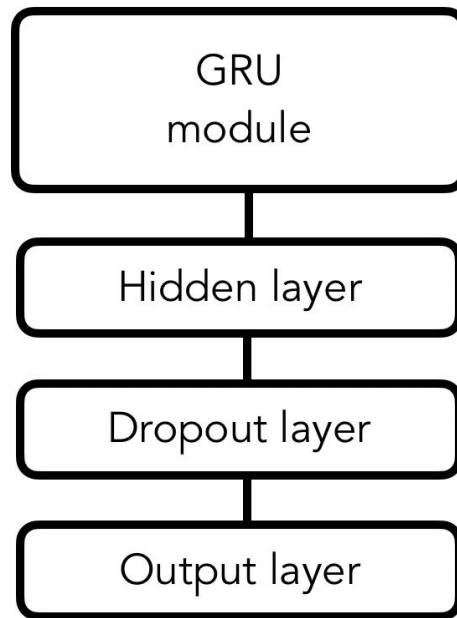


Figure 4. Single GRU model architectural diagram.

Stacked LSTM modules model

Our team developed this LSTM-based model in parallel with the observe 80/predict 30 model. This one is composed of an LSTM module whose output and final state are fed into a second LSTM module. The output of this second LSTM is input into a series of three dense layers. During the training phase, the model uses a dropout layer between the dense layers. During training, data was passed to the model in batches of 128 full length (110 frame) sequences.

The conceptualization here was to mimic the Observe 80/predict 30 model's ability to predict forecasting data over longer intervals. The previous models to this one predicted one new timestep for each batch of information sent as input. After training the model on the full computed dataset, it was validated using the full version of the validation set but with each agent time series providing 80 observations to the input (as opposed to 110 frames). The methodology here was born out of the 80/30 nature of observe/predict architecture. It should be noted that experiments were performed as well where unredacted observation series. The difference in testing metrics appeared to be negligible.

The model is recursively called 30 times to build the one observation output into a series of 30 forecasts. We experimented with predicting 2 (position_x, position_y) and 5 (position_x, position_y, MIN_DISTANCE_FRONT, MIN_DISTANCE_BACK, and NUM_NEIGHBORS) features to address challenges in the development. These discussion of these challenges will be expanded upon in later sections

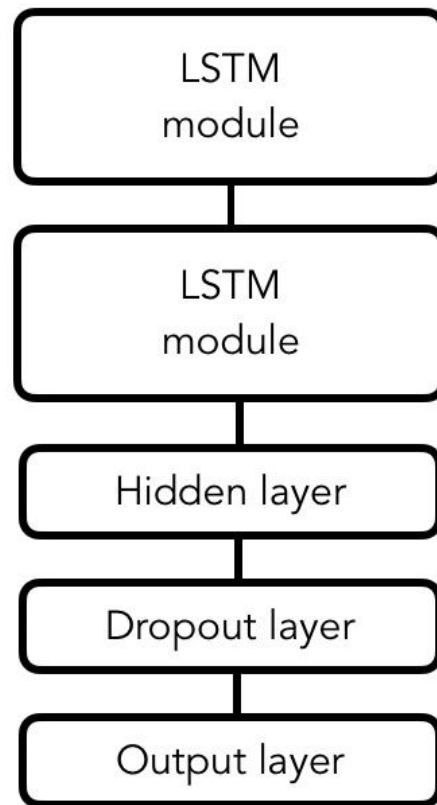


Figure 5. Stacked LSTM model architectural diagram.

Observe 80/predict 30 model

The observe/predict model has three components, a GRU, a GRU cell, and a set of dense layers. The observer processes the hidden layers, and returns a memory state over the entire observation sequence. The GRU cell takes the last observed values and the hidden state, as input state and input memory state. Over the entire prediction process, the cell uses the previous output states as input. At each prediction of the 30 prediction steps, the model calls the dense layers from the GRU cell's output state to produce a set of two values.

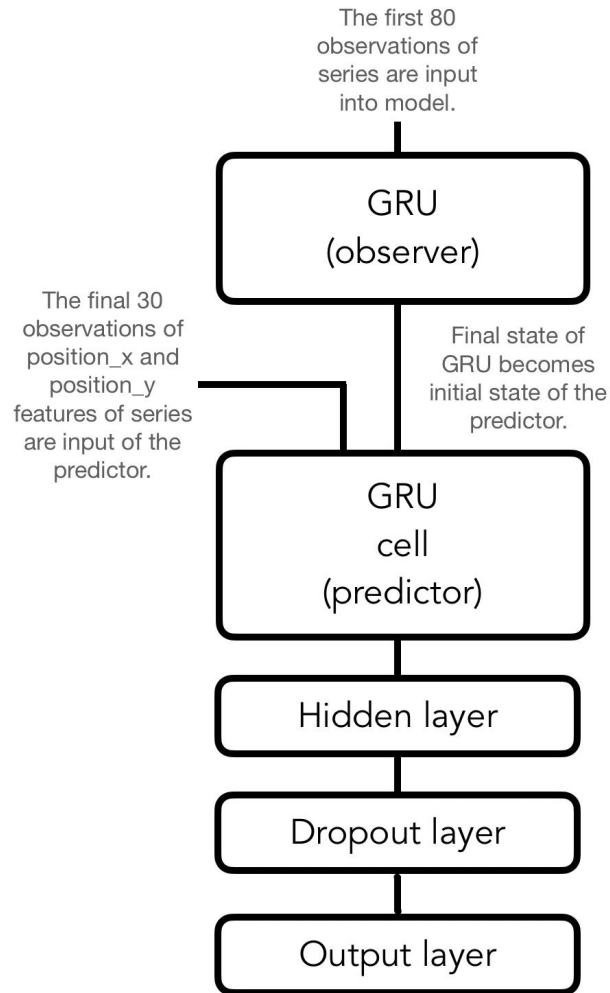


Figure 6. Observe 80 / predict 30 model architectural diagram.

Evaluation Metrics

Average displacement error, or ADE, is the average root mean squared error along a trajectory. It measures the difference between predicted and actual trajectories without exaggerating the effects of outliers.

MSE is used during our training process. It exaggerates the impact of large outliers during training so that the model can learn to minimize outlining inputs with reasonable speed.

Results

The observe 80/predict 30 model produces largely accurate predictions. We're observing a small degree of variability from between the ground truth and the training result as we move further from the last point of observation.

The Stacked LSTM performed comparably on metrics^[6] when attempting to predict only x and y positions. When the number of predicted features increased to 5 (which was necessary to attempt to implement the recursive forecasting strategy) its performance suffered in comparison.

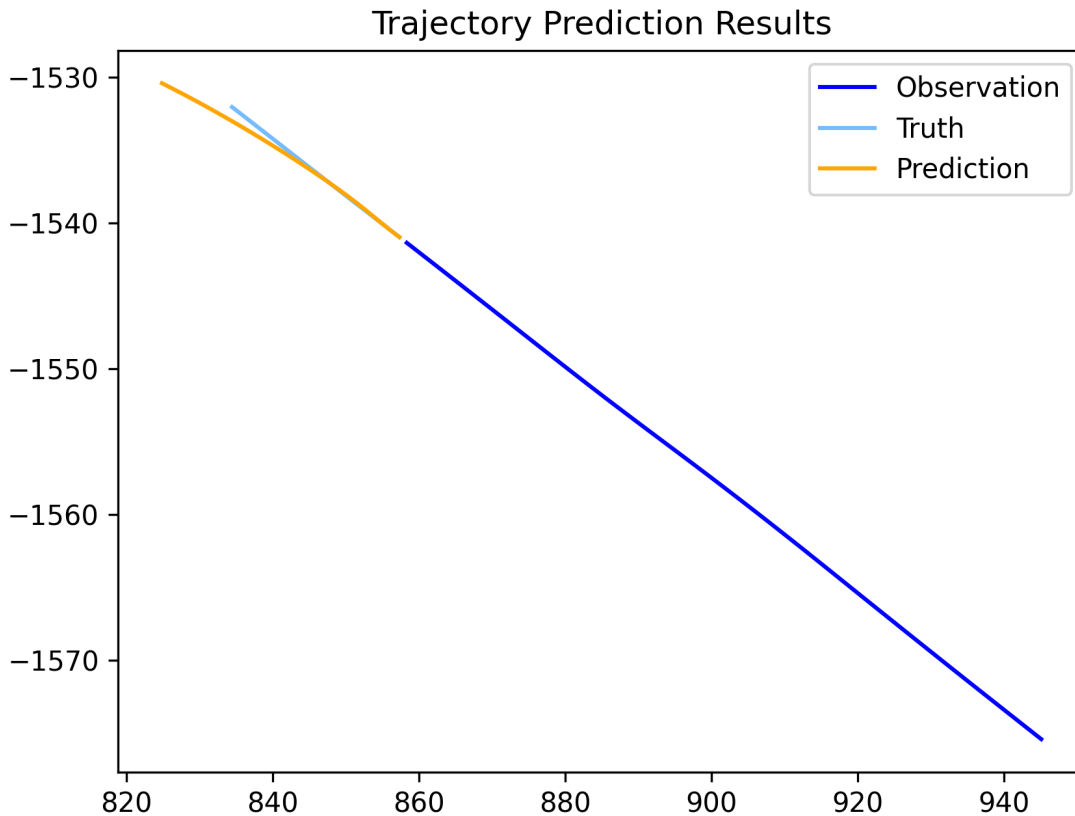


Figure 7. Trajectory prediction results plot for Observe 80/predict 30 model.

The following figure summarizes our validation metrics:

Model	MSE	ADE
Observe/Prediction	1.23	1.06
Stacked LSTM predicting 2 features (80 validation observations)	1.22	1.23
Stacked LSTM predicting 5 features (80 validation	255.05	22.33

observations)		
---------------	--	--

Figure 8. Metrics summary.

Discussion

In the observation/training model, since the next step is completely based on the results from the output state from the previous step, we worry that the error between the truth and the predicted positions will gradually build up. We did observe such errors in prediction results with particularly small unit size on particularly error-prone datapoints:

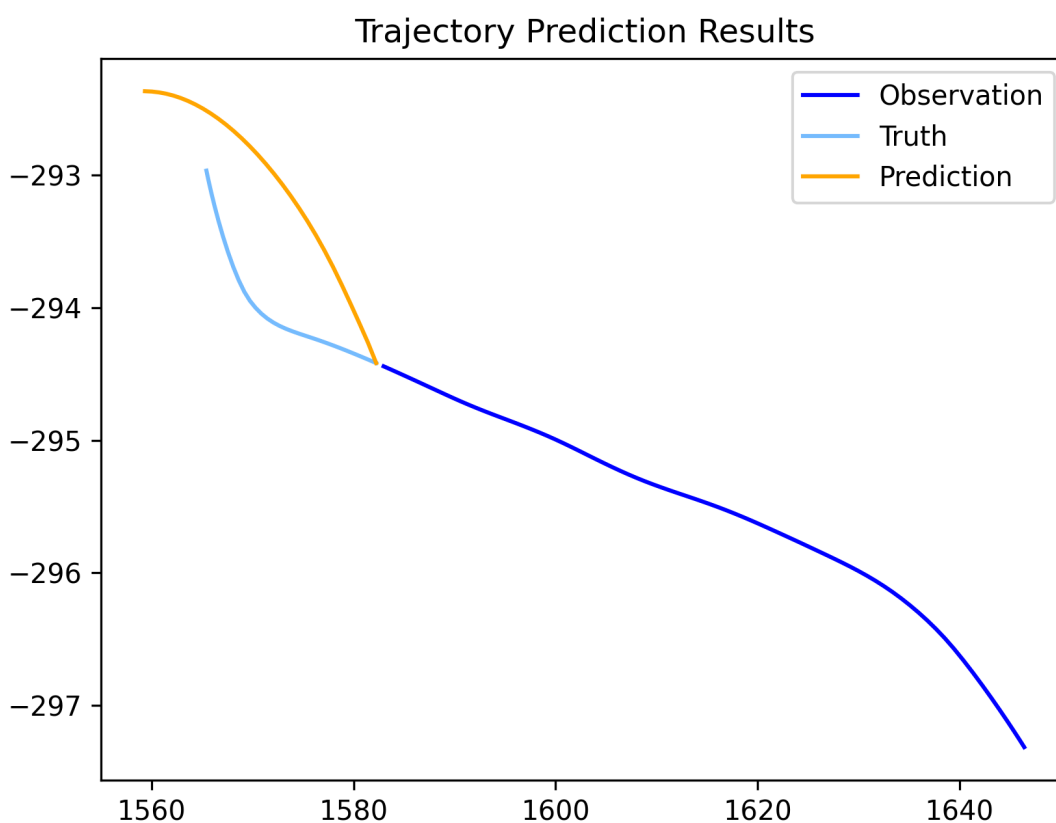


Figure 9.

We achieved higher ADE metrics than any models reported in the Argoverse 1 leaderboard (DCMS model at 1.4768), by a large amount. The observe and predict cadence creates a separation of concerns. The predictor gets the benefit of learning from the context of the observer while also only receiving 2 values for positional input into the cell while outputting 2

predictive values. The dense layers are also repeatedly called for each prediction step, together with the predictor GRU cell. This directly addresses the issues of recursive forecasting that the other models were hampered by the recursive prediction task.

We thought the prediction/observation model was the most reasonable, built it, and discovered its performance improvements. However, we are not sure which component or some nature of the dataset caused such performance improvement. To explain our performance, we performed the following experiments.

In the first experiment, we perform an ablation test on the recursive dense layers. We remove it from the recursive step. Instead, we use our GRU cells to recursively encode and output hidden information at a dimension of 128, and have the dense layers make final predictions at the end. We can see that the experimental setup performs better in ADE than our main model, and worse in MSE, indicating slightly worse performance for edge cases.

Experimental Setup	MSE	ADE
Recursive dense layers	1.23	1.06
Flat dense layer	1.28	1.00

If we significantly decrease our observation length, we can see that our accuracies slightly increase. However, the accuracy is still very high.

Experimental Setup	MSE	ADE
Observe 80 Predict 30	1.23	1.06
Observe 20 Predict 30	1.82	1.20

When we use only x,y positions to make forecasts, we understand that we could be training a neural network to perform regression without other contexts. But we noticed something very surprising:

Experimental Setup	MSE	ADE
With social features	1.23	1.06
Without social features	0.56	0.57

It turns out that the neural network trained without social features can perform significantly better in cases with nonlinear trajectories:

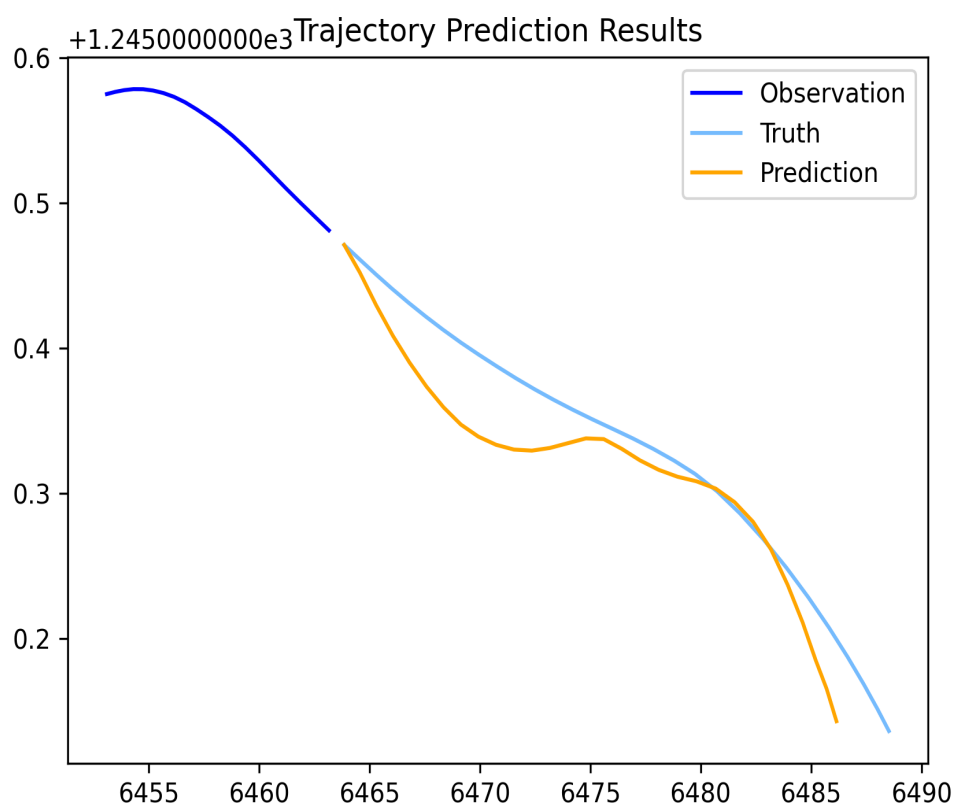


Figure 10.

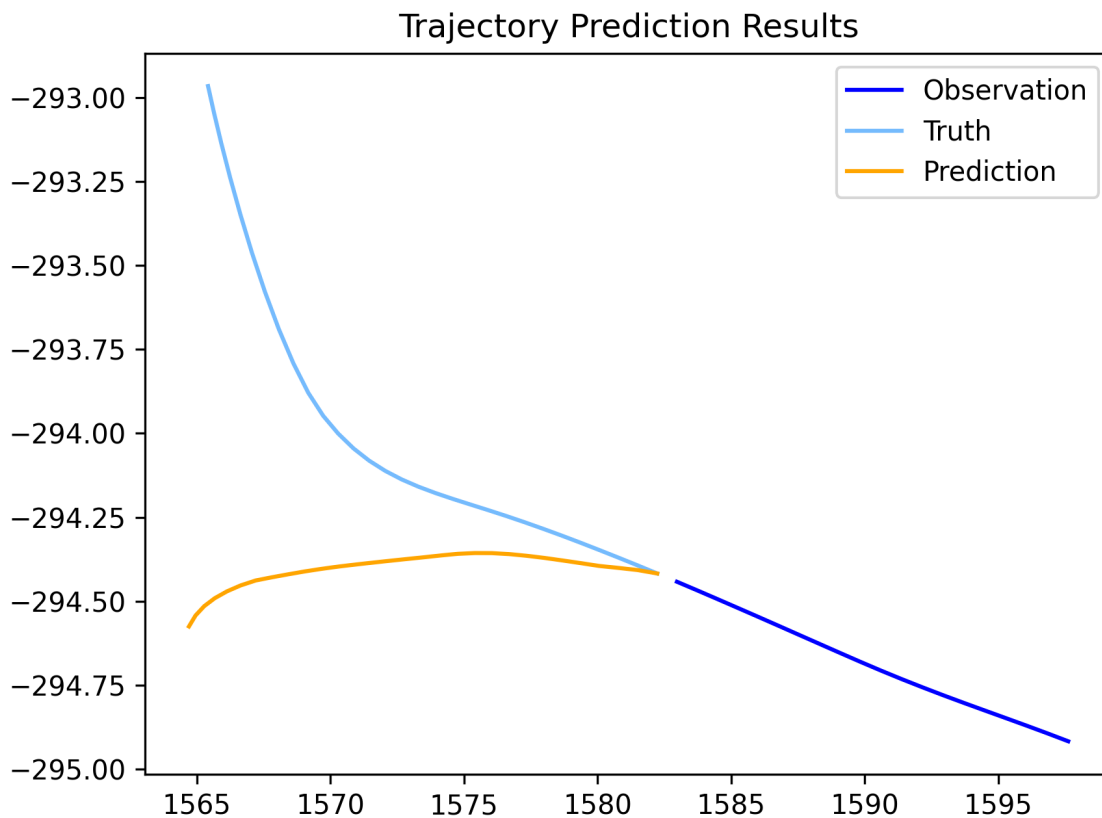


Figure 11.

In the meantime, it performs worse in other scenarios. The scenario above compares with Figure 9.

We believe that the network trained without social features performs better due to the large performance gain it obtains by performing tasks similar to linear regression in some cases but still failing in other cases.

Challenges

Argoverse 2 is a novel dataset and motion forecasting in self-driving is a novel topic for everyone in our team. Based on our prediction results (for instance, Figure 9) and how our models perform given the existing metrics, we believe that there is still a lot of future research work to develop better metrics, to systematically rule out large difference between observation and prediction.

The Stacked LSTM model was challenging to work into a true motion forecasting task using our method of recursive calls predicting the next step from the previous. For example, the architecture dictates that the input must include all 5 of the feature variables used to train the model. If the model is tasked to only predict 2 positional features then the process halts when attempting to feed that output back into the input recursively because the computed social features are not present. We experimented with concatenating values of 0 to represent the social features to the output before the next recursive step. Value imputation may also be an option here since in essence the social features are “missing” for the recursive input. Both options are potentially dangerous in this context.

Safety is paramount in applications of autonomous functionality. The two methods mentioned above cannot guarantee accuracy of prediction because the social feature values are either arbitrary or potentially inaccurate. The results show that it is possible to take the strategy of predicting all 5 features to satisfy the recursive step but as evidenced the metric accuracy suffers significantly.

Reflection & Future Work

Forecasting, by essence, is a very difficult task. For humans to forecast motion, we typically rely on the simple fact that an object will continue moving at the same speed as we see it. Based on our observations, recurrent networks perform well in learning this rule (for instance, figure 7). Therefore we believe that most of the challenge in this domain would be focused on cases such as Figure 9. However, it is worth debating whether there is any pattern based on nonlinear trajectories at all. In other words, we need to both explain how humans would perform such forecasting and in this context explain neural networks would be able to learn realistic patterns and make predictions.

A lesson we learned is to measure the amount of data and training time per batch early, so we could more accurately reason about compute time and resources. We also noticed that both GCP and Oscar are down near the end of the semester; it would be helpful to find and use a reliable and free alternative computational infrastructure at the very early stage of the project.

We also gained valuable experience in designing training pipelines. We found that it is extremely important to log and visualize metrics, and to save trained weights in an organized fashion. In addition, we also found that it is helpful to serve data as a generator, either with a Dataset class or a Python generator function.

Additionally, we were initially hoping to get to the point where we could adversarially attack^[7] our models, but we ran out of time. Pivoting away from this helped us focus on getting reportable results from several of the models we had in development at the time, which was a good choice.

Despite the large improvement in ADE, we believe that there is still a lot of room for improvements, especially to improve accuracy for cases with non-linear trajectories such as shown in Figure 9.

It's also interesting to consider the application of observe and predict strategies to other forms of time series data such as stock prices, audio, and video.

Citation

1. Argoverse 2 Dataset:
<https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/4734ba6f3de83d861c3176a6273cac6d-Paper-round2.pdf>
2. Argo AI, <https://www.argoverse.org/>
3. Multimodal Motion Prediction with Stacked Transformers Paper:
<https://arxiv.org/pdf/2103.11624.pdf>
4. How to Develop LSTM Models for Time Series Forecasting
<https://machinelearningmastery.com/how-to-develop-lstm-models-for-time-series-forecasting/>
5. How to Develop Multilayer Perceptron Models for Time Series Forecasting
<https://machinelearningmastery.com/how-to-develop-multilayer-perceptron-models-for-time-series-forecasting/>
6. Argoverse 1 leaderboard <https://eval.ai/challenge/454/leaderboard/1279>
7. https://openaccess.thecvf.com/content/ICCV2021/papers/Li_Fooling_LiDAR_Perception_via_Adversarial_Trajectory_Perturbation_ICCV_2021_paper.pdf