



HAL
open science

INTELLinext: A Fully Integrated LSTM and HMM-Based Solution for Next-App Prediction With Intel SUR SDK Data Collection

Cyril Gorlla, Jared Thach, Hiroki Hoshida

► **To cite this version:**

Cyril Gorlla, Jared Thach, Hiroki Hoshida. INTELLinext: A Fully Integrated LSTM and HMM-Based Solution for Next-App Prediction With Intel SUR SDK Data Collection. Halçioğlu Data Science Institute Capstone Showcase, Apr 2022, La Jolla, United States. hal-03825175

HAL Id: hal-03825175

<https://hal.science/hal-03825175v1>

Submitted on 22 Oct 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

INTELLinext: A Fully Integrated LSTM and HMM-Based Solution for Next-App Prediction With Intel SUR SDK Data Collection

Cyril Gorlla
Jared Thach
Hiroki Hoshida
cyril.m.gorlla@jacobs.ucsd.edu
j1thach@ucsd.edu
hhoshida@ucsd.edu
Halicioğlu Data Science Institute
University of California San Diego
La Jolla, California, USA

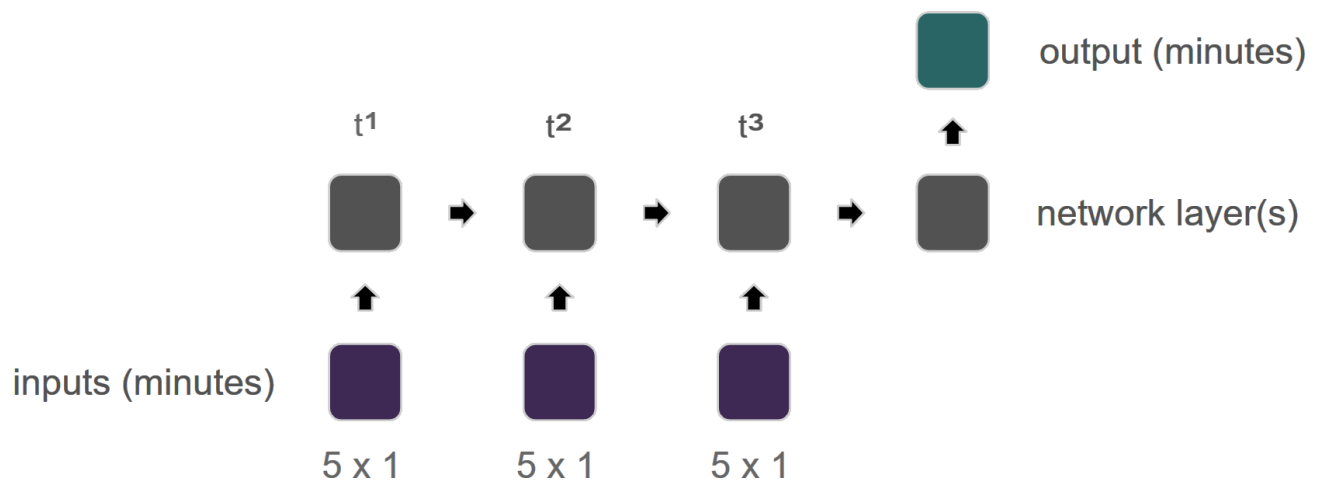


Figure 1. Diagram of LSTM for usage duration prediction.

CCS Concepts: • Computing methodologies → Machine learning; Neural networks; Markov decision processes; Active learning settings.

Keywords: app prediction, lstm, hidden markov model

1 Abstract

As the power of modern computing devices increases, so too do user expectations for them. Despite advancements in

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

HDSI Intel Capstone '21, Jan 3–March 19, 2022, San Diego, CA

© 2022 Copyright held by the owner/author(s).

technology, computer users are often faced with the dreaded spinning icon waiting for an application to load. Building upon our previous work developing data collectors with the Intel System Usage Reporting (SUR) SDK, we introduce INTELLinext, a comprehensive solution for next-app prediction for application preload to improve perceived system fluidity. We develop a Hidden Markov Model (HMM) for prediction of the k most likely next apps, achieving an accuracy of 70% when $k = 3$. We then implement a long short-term memory (LSTM) model to predict the total duration that applications will be used. After hyperparameter optimization leading to an optimal lookback value of 5 previous applications, we are able to predict the usage time of a given application with a mean absolute error of 45 seconds. Our work constitutes a promising comprehensive application preload solution with data collection based on the Intel SUR SDK and prediction with machine learning.

2 Introduction

Users often face a myriad of minuscule slow-downs when navigating a computer. These slow-downs largely take the form of application loading times and we aimed to first collect user data by building custom data collectors in our previous work. The collected data can then be used to study user patterns and to extract insights via machine learning practices and models such as Hidden Markov Models (HMM) and Long Short-Term Memory (LSTM) Models. By leveraging these models, the inconvenience of accumulated application start-up times can be greatly reduced.

Despite computer processor speeds increasing year after year, there are still high usage programs and processes that interrupt workflows in our daily lives. Whether that interruption be lag when opening Zoom to join a meeting, or waiting for Google Chrome to open a link embedded in a Word document, these micro (or in some cases, major) stutters can cause a process meant to be relatively smooth to be a large source of frustration in the daily lives of an end user. One key point, however, is that these processes are actually often “scheduled” in a sense. That is, most people often repeat the same or similar tasks every day, and have a set routine. One example could be an office worker who opens Microsoft Excel almost every time they turn on their laptop at the office. These sorts of patterns can be studied, and by using machine learning, we can create predetermined schedules for users, allowing us to preload apps, or have processes ready before the user needs them. In other words, by analyzing user behavior, we can load an application the user would likely use next in advance, reducing wait and loading times.

By using Intel’s System Usage Reporting library and the accompanying XLS SDK, we created “monitors” of user and computer activity, known as collectors or Input Libraries, and create activity logs which we can then analyze to provide preloading solutions for the user. In our previous work, we have covered the development of these Input Libraries. We now briefly restate their descriptions.

In order for us to create our own Input Libraries, the XLS SDK provides a wide range of examples which can be used either as references or templates. Throughout the first ten weeks of working with Intel, we developed four different Input Libraries, each using a different template and measuring different categories of inputs from the user and computer. The first, the `mouse_input` Input Library, keeps a log of the cursor coordinates as the user moves the mouse around. Second, the `user_waiting` Input Library, keeps a timer based log of the cursor icon as the user uses the computer. The third, a `foreground_window` Input Library, creates a log entry whenever the foreground window (the window in front of all other windows) changes whether it be automatically (such as a notification pop up) or by user input (clicking the taskbar). Finally, the fourth Input Library is the

`desktop_mapper`, which, when triggered by a change in the foreground window, maps all the windows on the desktop in z-order and stores pertinent information about each window e.g. position and size. Each of these Input Libraries are coded differently in fundamental ways, and measure changes in different ways as well. By using the data provided by Libraries like these, we can determine preloading schedules for the individual user.

3 Methods

3.1 Data Collection & Preliminary Analysis

With our Input Libraries fully functional, they were continuously run on two computers (identified by the IDs LAPTOP and DESKTOP) over the course of three months, from November 2021 to February 2022. The users were advised to continue operating the machine according to their regular schedules to ensure that the following analysis would best generate insights on naturally occurring user patterns. In total, the Input Libraries collected over 160,000 rows of raw data in 76 database (.db) files from the two machines. Thanks to our extensive data validation process we built in the first half of the project, the data cleanup needed was minimal. These processes include functions to check for appropriate data types, improbable values (such as negative mouse X and Y coordinates), and null values. As a result, the manual work required mainly consisted of column renaming and other simple DataFrame manipulations. After extracting just the foreground process information, we were left with two datasets with 5,241 and 10,113 rows respectively.

We found that each computer had separate periods of high and low usage, with general boundaries around Fall Quarter 2021, Winter Break 2021/2022, and Winter Quarter 2022. Because of the higher consistency of data of some periods depending on the machine, we used data before December 10, 2021 for the laptop and we used data after January 7, 2022 for the desktop.

Some trends in the data can be seen in the exploratory plots below. For instance, Chrome was the most common application in the dataset followed by Windows Explorer for both datasets, but the trends differ after those. We found that each computer had separate periods of high and low usage, with general boundaries around Fall Quarter 2021, Winter Break 2021/2022, and Winter Quarter 2022. Because of the higher consistency of data of some periods depending on the machine, we used data before December 10, 2021 for the laptop and we used data after January 7, 2022 for the desktop. Some trends in the data can be seen in the exploratory plots below. For instance, Chrome was the most common application in the dataset followed by Windows Explorer for both datasets, but the trends differ after those.

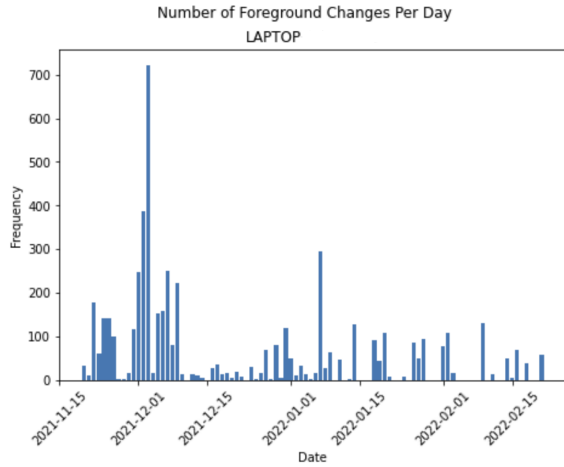


Figure 2. Laptop foreground changes per day

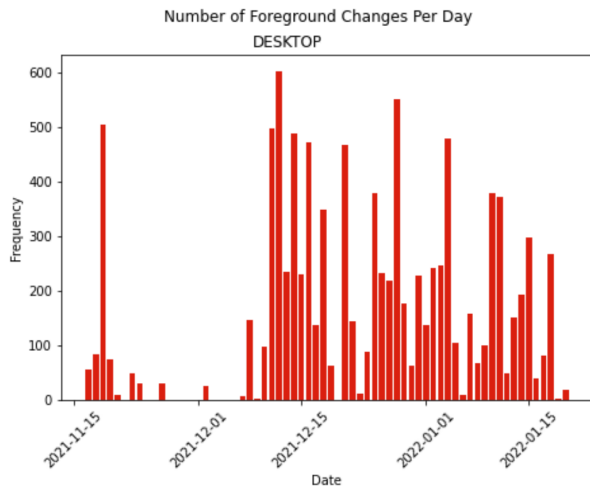


Figure 3. Desktop foreground changes per day

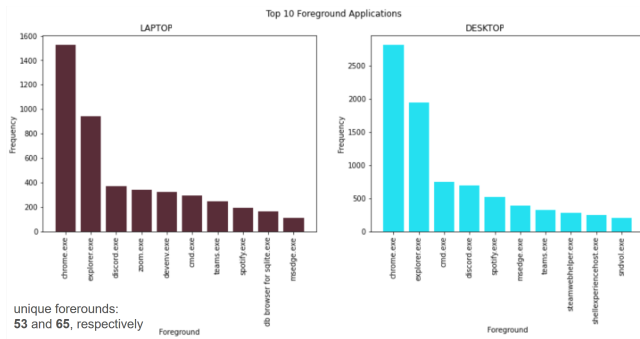


Figure 4. Top 10 foregrounds for each machine

3.2 Model Building & Data Science

We addressed two predictive tasks using the data collected. For our first task, we built models to predict the next application used based on previous user activity. For our second task, we built a model to predict the duration the next application would be used based on previous user activity.

3.2.1 Task 1. Approach 1. Next-App Prediction: Hidden Markov Model

The Hidden Markov Model (HMM) is a particular machine learning approach that ingests sequences of time-related events in order to predict future events. For our purposes, we manually implemented a First Order HMM (adopting a similar model class as the package, scikit-learn) which looks at a single previous application in order to predict the single next application. A HMM’s “order” refers to the amount of previous applications, or “look-back”, the HMM will use in order to generate a single next prediction, and therefore, our First Order HMM will use a single previous event to predict a single next event. This is equivalent to simply computing conditional probabilities of each unique 2 sequence event; for a given previous event A, the probability of the next event B is determined by:

$$P(B | A) = \frac{P(B \cap A)}{P(A)} \tag{1}$$

These probabilistic values are calculated upon model training and stored into a posterior matrix instance variable. These probabilities are then recalled upon prediction.

Given an input foreground of notepad++.exe, for example, the probability of explorer.exe being the next foreground is approximately 43% as shown in Figure 4. One special implementation of our First Order HMM is the custom predict function which has an optional parameter of n_foregrounds which specifies the number of foregrounds to return for a single input, based on the foregrounds with the highest probabilities. When viewing Figure 4 yet again, a predict with n_foregrounds = 3 will return the list [‘explorer.exe’, ‘mmc.exe’, ‘chrome.exe’]. Accuracy for a single observation is therefore calculated by determining whether or not the true foreground application exists within the list of predicted foreground applications. Using n_foregrounds = 3 for our testing dataset, we achieve a final accuracy of 70.05%.

3.2.2 Task 1. Approach 2. Next-App Prediction: Long Short-Term Memory

Though our Hidden Markov Model implementation was fairly successful, we wanted to approach the task from a different angle, so we built a Long Short-Term Memory Recurrent Neural Network (LSTM RNN) model as well. Because of the higher complexity of LSTMs compared to HMMs, we developed our model using Keras, a deep learning API built on Python that allows for a streamlined pipeline of model

creation. By using Keras, we could manipulate various values quickly and easily, allowing us to test many variations of the models in a short period of time, without having to rewrite any significant part of the code. Keras also contains many preprocessing and testing functions that allow us to prepare our data and test the model accuracy easily. We decided on a RNN model because RNNs are a type of neural network that uses previous inputs to make new predictions. This was perfect for our case, as our data is time based, and we were trying to make future predictions based on past actions. Furthermore, out of the many implementations of RNNs, we decided it was best to use a Long Short-Term Memory (LSTM) model. This is because LSTM models are more tailored to situations where the durations between events vary. In our case, the durations between each foreground application switch were different, meaning that the timestamps of each point of data were spread unevenly.

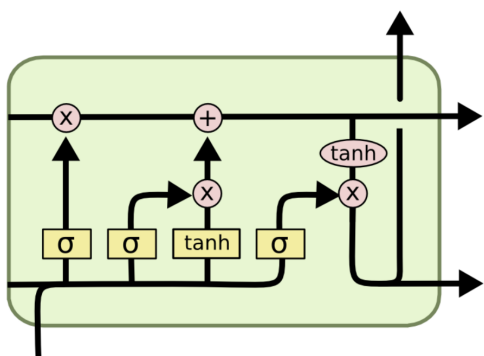


Figure 5. LSTM cell (Source: Christopher Olah)

In order to use our data to create the model, we first needed to process it into the correct format so Keras could read it in. First, we needed to select the features to use for our model. In other words, the model would use these properties to make the prediction. We used the time of foreground change, the previous application used, and the time the user spent on the previous application right before the foreground window change. Even with the features selected, however, we had to do further preprocessing to make the data usable. First, we extracted just the hour out of the timestamp. This is because we wanted to see if app usage depended on time of day, and removing non-repeating elements (such as month and day) and removing elements too minute (such as minutes and seconds) seemed appropriate. Second, we one-hot encoded the previous application used. As we found that some of the applications found in the dataset only appeared once or twice (installers, for example), we found the top ten most used applications, and renamed the rest to “Other”. Though initially we did not do this, after multiple rounds of testing we found that removing the lower outliers improved our model accuracy tremendously. Finally, we created the output

column by shifting the application name column by one. After preprocessing, we split the data into a 70:30 ratio, in which 70% of the data was used for training, and 30% was used for testing. The training dataset was used to create the LSTM model in Keras. Our final model consisted of four layers; an LSTM layer, a dropout layer, another LSTM layer, and a dense layer with a softmax activation function. Thanks to the flexibility and ease of use of Keras layers, we were able to experiment with various configurations, and we found that using the ADAM optimizer, the categorical cross-entropy loss function, and 100 epochs created a fairly accurate model of our training set.

In order to use our data to create the model, we first needed to process it into the correct format so Keras could read it in. First, we needed to select the features to use for our model. In other words, the model would use these properties to make the prediction. We used the time of foreground change, the previous application used, and the time the user spent on the previous application right before the foreground window change. Even with the features selected, however, we had to do further preprocessing to make the data usable. First, we extracted just the hour out of the timestamp. This is because we wanted to see if app usage depended on time of day, and removing non-repeating elements (such as month and day) and removing elements too minute (such as minutes and seconds) seemed appropriate. Second, we one-hot encoded the previous application used. As we found that some of the applications found in the dataset only appeared once or twice (installers, for example), we found the top ten most used applications, and renamed the rest to “Other”. Though initially we did not do this, after multiple rounds of testing we found that removing the lower outliers improved our model accuracy tremendously. Finally, we created the output column by shifting the application name column by one. After preprocessing, we split the data into a 70:30 ratio, in which 70% of the data was used for training, and 30% was used for testing. The training dataset was used to create the LSTM model in Keras. Our final model consisted of four layers; an LSTM layer, a dropout layer, another LSTM layer, and a dense layer with a softmax activation function. Thanks to the flexibility and ease of use of Keras layers, we were able to experiment with various configurations, and we found that using the ADAM optimizer, the categorical cross-entropy loss function, and 100 epochs created a fairly accurate model of our training set. Similarly to our HMM, our LSTM returned a list of the top four most likely applications to follow the current foreground application, ordered in levels of confidence. Accuracy was calculated with the same function (where a single observation was labeled accurate if the true future foreground application appeared in the list of predictions). Our LSTM model had a test accuracy of 68.60%, a value similar to our HMM accuracy.

3.2.3 Task 2. Next-App Duration Prediction: Long Short-Term Memory

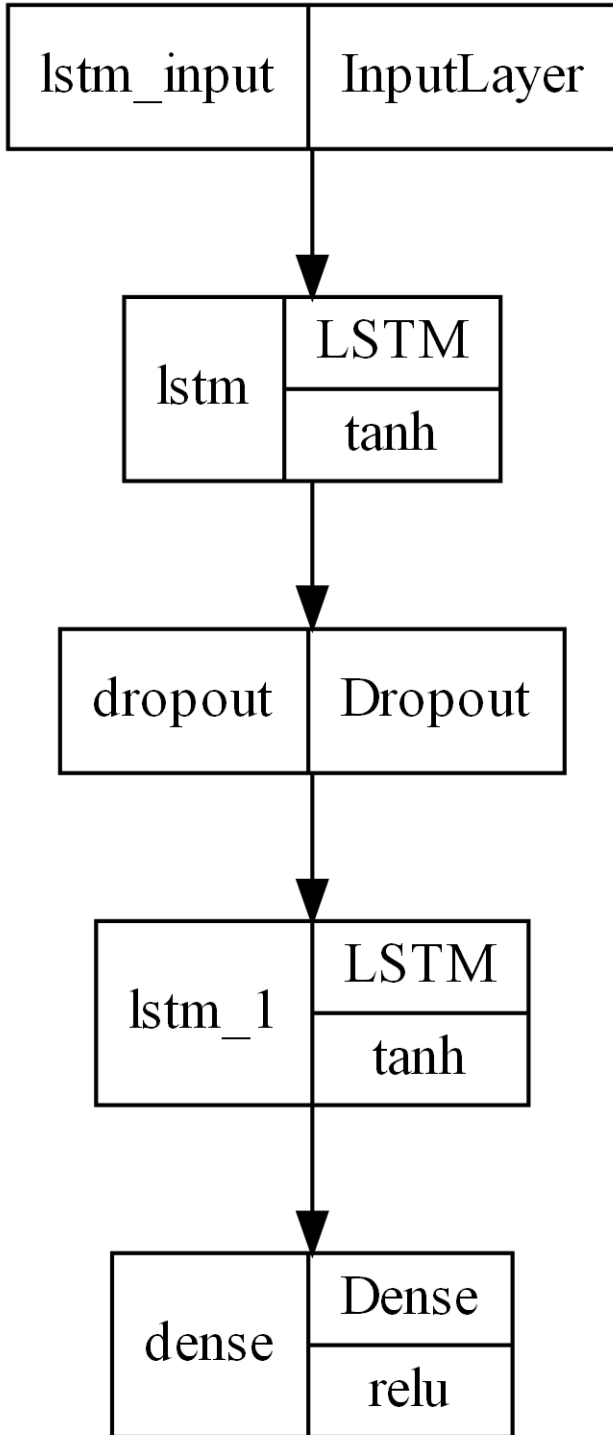


Figure 6. LSTM duration prediction model

Because of the flexibility of the LSTM model, we decided to focus on a LSTM solution for task 2, predicting next-application durations, as well. Our task 1 LSTM model used a “look-back” value of one previous foreground application in order to predict one future foreground application, where a “look-back” is defined as the number of previous events a single input will use in order to generate the next output. In order to raise accuracy, our task 2 LSTM used a look-back value of five. In other words, the model uses the previous five data points to predict the next. The task 2 model architecture is similar to our task 1 model, with the four layers in the same order. However, we used a MAE loss function, ReLU activation function, and 25 epochs instead. After training, our model made predictions with a mean absolute error of 0.74 minutes. With our testing dataset’s mean foreground duration of 1.73 minutes, on average our model predicted foreground durations accurate within approximately 44 seconds.

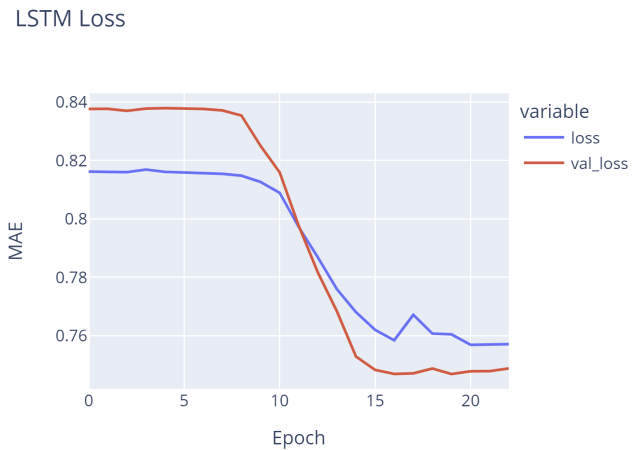


Figure 7. LSTM duration prediction model training and validation loss)

4 Results

Our project addressed two predictive tasks: (1) next-app prediction and (2) next-app duration prediction. Task 1 was explored via two machine learning models with our First Order Hidden Markov Model yielding an accuracy of 70.05% using a prediction list of three and our Long Short-Term Memory Model yielding an accuracy of 68.60% using a prediction list of four. Task 2 was explored through a single Long Short-Term Memory Model with an average error rate of 42.77%, or 44.4 seconds.

5 Conclusion

By leveraging C programming to custom code data collectors, we were able to efficiently collect user data while minimizing extraneous information gathering. Additionally, primary

data sourcing was elucidated as an essential aspect of machine learning practices.

Although we have successfully solved our two project tasks of (1) next-app prediction and (2) next-app duration prediction, there exists improvements in several realms. Our original HMM was based on a First Order HMM, and therefore, only used single, previous applications to predict future applications. By spending more time working towards a Second, Third, or higher Order HMM, we can hope to achieve greater model performance. Additionally during our data collection phase, there was heavy overlap of distinct time periods, such as Fall Quarter 2021, Winter Break, and Winter

Quarter 2022 (based on a college quarter system). In the future, we may aim to clearly demarcate data collection boundaries to ensure quality data that is independent from each distinct time period. With the aforementioned model performances for our classification-based next-app prediction task and regression-based next-app duration prediction task, we hope to implement our data science work for other relevant applications. These other applications include edge device machine learning, cloud machine learning, and repurposing of obsolete hardware.

Acknowledgments

Jamel Tayeb and Bijan Arbab, Intel; Intel DCA Team