




Artificial Neural Networks Applied to Taxi Destination Prediction

Alexandre de Brébisson¹ , Étienne Simon² , Alex Auvolat³ ,
Pascal Vincent^{1,4}, and Yoshua Bengio^{1,4}.


¹ MILA lab, University of Montréal,
alexandre.de.brebisson@umontreal.ca,
vincentp@iro.umontreal.ca

² ENS Cachan,
esimon@esimon.eu

³ ENS Paris,
alex.auvolat@ens.fr

⁴ CIFAR.

Abstract. We describe our first-place solution to the ECML/PKDD discovery challenge on taxi destination prediction. The task consisted in predicting the destination of a taxi based on the beginning of its trajectory, represented as a variable-length sequence of GPS points, and diverse associated meta-information, such as the departure time, the driver id and client information. Contrary to most published competitor approaches, we used an almost fully automated approach based on neural networks and we ranked first out of 381 teams. The architectures we tried use multi-layer perceptrons, bidirectional recurrent neural networks and models inspired from recently introduced memory networks. Our approach could easily be adapted to other applications in which the goal is to predict a fixed-length output from a variable-length sequence.

 Random order, this does not reflect the weights of contributions.

1 Introduction

The taxi destination prediction challenge was organized by the 2015 ECML/PKDD conference⁵ and proposed as a Kaggle competition⁶. It consisted in predicting the destinations (latitude and longitude) of taxi trips based on initial partial trajectories (which we call *prefixes*) and some meta-information associated to each ride. Such prediction models could help to dispatch taxis more efficiently.

The dataset is composed of all the complete trajectories of 442 taxis running in the city of Porto (Portugal) for a complete year (from 2013-07-01 to 2014-06-30). The training dataset contains 1.7 million datapoints, each one representing a complete taxi ride and being composed of the following attributes⁷:

- the complete taxi ride: a sequence of GPS positions (latitude and longitude) measured every 15 seconds. The last position represents the destination and different trajectories have different GPS sequence lengths.
- metadata associated to the taxi ride:
 - if the client called the taxi by phone, then we have a client ID. If the client called the taxi at a taxi stand, then we have a taxi stand ID. Otherwise we have no client identification,
 - the taxi ID,
 - the time of the beginning of the ride (unix timestamp).

In the competition setup, the testing dataset is composed of 320 partial trajectories, which were created from five snapshots taken at different timestamps. This testing dataset is actually divided in two subsets of equal size: the public and private test sets. The public set was used through the competition to compare models while the private set was only used at the end of the competition for the final leaderboard.

Our approach uses very little hand-engineering compared to those published by other competitors. It is almost fully automated and based on artificial neural networks. Section 2 introduces our winning model, which is based on a variant of a multi-layer perceptron (MLP) architecture. Section 3 describes more sophisticated alternative architecture that we also tried. Although they did not perform as well as our simpler winning model for this particular task, we believe that they can provide further insight on how to apply neural networks to similar tasks. Section 4 and Section 5 compares and analyses our various models quantitatively and qualitatively on both the competition testing set and a bigger custom testing set.

⁵ <http://www.geolink.pt/ecmlpkdd2015-challenge/>

⁶ <https://www.kaggle.com/c/pkdd-15-predict-taxi-service-trajectory-i>

⁷ The exact list of attributes for each trajectory can be found here: <https://www.kaggle.com/c/pkdd-15-predict-taxi-service-trajectory-i/data>

2 The Winning Approach

2.1 Data Distribution

Our task is to predict the destination of a taxi given a prefix of its trajectory. As the dataset is composed of full trajectories, we have to generate trajectory prefixes by cutting the trajectories in the right way. The provided training dataset is composed of more than 1.7 million complete trajectories, which gives 83 480 696 possible prefixes. The distribution of the training prefixes should be as close as possible as that of the provided testing dataset on which we were eventually evaluated. This test set was selected by taking five snapshots of the taxi network activity at various dates and times. This means that the probability that a trajectory appears in the test set is proportional to its length and that, for each entire testing trajectory, all its possible prefixes had an equal probability of being selected in the test set. Therefore, generating a training set with all the possible prefixes of all the complete trajectories of the original training set provides us with a training set which has the same distribution over prefixes (and whole trajectories) as the test set.

2.2 MLP Architecture

A Multi-Layer Perceptron (MLP) is a neural net in which each neuron of a given layer is connected to all the neurons of the next layer, without any cycle. It takes as input fixed-size vectors and processes them through one or several hidden layers that compute higher level representations of the input. Finally the output layer returns the prediction for the corresponding inputs. In our case, the input layer receives a representation of the taxi's prefix with associated metadata and the output layer predicts the destination of the taxi (latitude and longitude). We used standard hidden layers consisting of a matrix multiplication followed by a bias and a nonlinearity. The nonlinearity we chose to use is the Rectifier Linear Unit (ReLU) [1], which simply computes $\max(0, x)$. Compared to traditional sigmoid-shaped activation functions, the ReLU limits the gradient vanishing problem as its derivative is always one when x is positive. For our winning approach, we used a single hidden layer of 500 ReLU neurons.

2.3 Input Layer

One of the first problems we encountered was that the trajectory prefixes are varying-length sequences of GPS points, which is incompatible with the fixed-size input of the MLP. To circumvent (temporarily, see Section 3.1) this limitation, we chose to consider only the first k points and last k points of the trajectory prefix, which gives us a total of $2k$ points, or $4k$ numerical values for our input vector. For the winning model we took $k = 5$. These GPS points are standardized (zero-mean, unit-variance).

To deal with the discrete meta-data, consisting of client ID, taxi ID, date and time information, we learn embeddings jointly with the model for each of

these information. This is inspired by neural language modeling approaches [2], in which each word is mapped to a vector space of fixed size (the vector is called the *embedding* of the word). The table of embeddings for the words is included in the model parameters that we learn and behaves as a regular parameter matrix: the embeddings are first randomly initialized and are then modified by the training algorithm like the other parameters. In our case, instead of words, we have metadata values. More precisely, we have one embedding table for each metadata with one row for each possible value of the metadata. For the date and time, we decided to create higher-level variables that better describe human activity: quarters of hour, day of the week, week of the year (one embedding table is learnt for each of them). These embeddings are then simply concatenated to the $4k$ GPS positions to form the input vector of the MLP. The complete list of embeddings used in the winning model is given in Table 1.

Table 1: Metadata values and associated embedding size.

Metadata	Number of possible values	Embedding size
Client ID	57106	10
Taxi ID	448	10
Stand ID	64	10
Quarter hour of the day	96	10
Day of the week	7	10
Week of the year	52	10

2.4 Destination Clustering and Output Layer

As the destination we aim to predict is composed of two scalar values (latitude and longitude), it is natural to have two output neurons. However, we found that it was difficult to train such a simple model because it does not take into account any prior information on the distribution of the data. To tackle this issue, we integrate prior knowledge of the destinations directly in the architecture of our model: instead of predicting directly the destination position, we use a predefined set $(c_i)_{1 \leq i \leq C}$ of a few thousand destination cluster centers and a hidden layer that associates a scalar value $(p_i)_i$ (similar to a probability) to each of these clusters. As the network must output a single destination position, for our output prediction \hat{y} , we compute a weighted average of the predefined destination cluster centers:

$$\hat{y} = \sum_{i=1}^C p_i c_i.$$

Note that this operation is equivalent to a simple linear output layer whose weight matrix would be initialized as our cluster centers and kept fixed during

training. The hidden values $(p_i)_i$ must sum to one so that \hat{y} corresponds to a centroid calculation and thus we compute them using a softmax layer:

$$p_i = \frac{\exp(e_i)}{\sum_{j=1}^C \exp(e_j)},$$

where $(e_j)_j$ are the activations of the previous layer.

The clusters $(c_i)_i$ were calculated with a mean-shift clustering algorithm on the destinations of all the training trajectories, returning a set of $C = 3392$ clusters. Our final MLP architecture is represented in Figure 1.

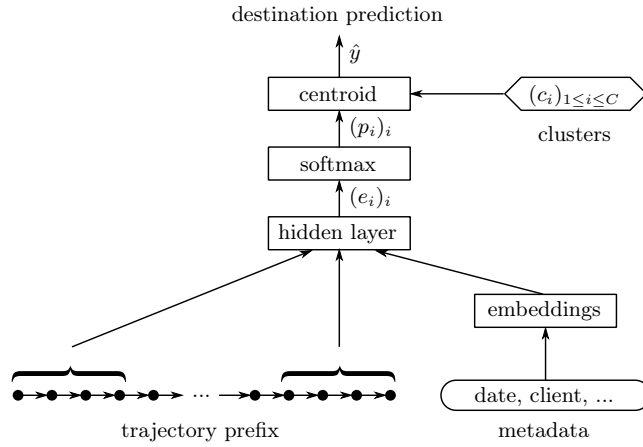


Fig. 1: Architecture of the winning model.

2.5 Cost Computation and Training Algorithm

The evaluation cost of the competition is the mean Haversine distance, which is defined as follows (λ_x is the longitude of point x , ϕ_x is its latitude, and R is the radius of the Earth):

$$d_{\text{haversine}}(x, y) = 2R \arctan \left(\sqrt{\frac{a(x, y)}{a(x, y) - 1}} \right),$$

where $a(x, y)$ is defined as:

$$a(x, y) = \sin^2 \left(\frac{\phi_y - \phi_x}{2} \right) + \cos(\phi_x) \cos(\phi_y) \sin^2 \left(\frac{\lambda_y - \lambda_x}{2} \right).$$

Our models did not learn very well when trained directly on the Haversine distance function and thus, we used the simpler equirectangular distance instead,

which is a very good approximation at the scale of the city of Porto:

$$d_{\text{equirectangular}}(x, y) = R \sqrt{\left((\lambda_y - \lambda_x) \cos\left(\frac{\phi_y - \phi_x}{2}\right) \right)^2 + (\phi_y - \phi_x)^2}.$$

We used stochastic gradient descent (SGD) with momentum to minimise the mean equirectangular distance between our predictions and the actual destination points. We set a fixed learning rate of 0.01, a momentum of 0.9 and a batch size of 200.

3 Alternative Approaches

The models that we are going to present in this section did not perform as well for our specific destination task on the competition test set but we believe that they can provide interesting insights for other problems involving fixed-length outputs and variable-length inputs.

3.1 Recurrent Neural Networks

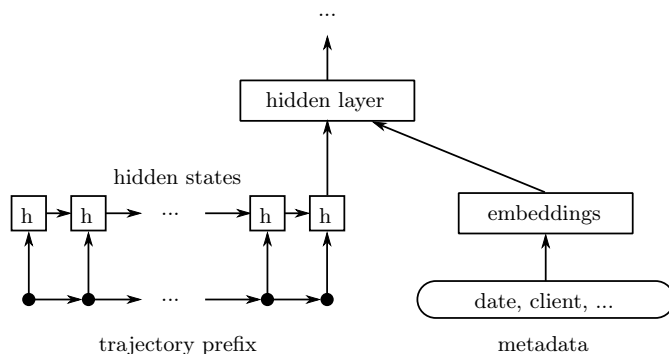


Fig. 2: RNN architecture, in which the GPS points are read one by one in the forward direction of the prefix. (The rest of the architecture is the same as before and is omitted for clarity.)

As stated previously, a MLP is constrained by its fixed-length input, which prevents us from fully exploiting the entire trajectory prefix. Therefore we naturally considered recurrent neural net (RNN) architectures, which can read all the GPS points one by one, updating a fixed-length internal state with the same transition matrix at each time step. The last internal state of the RNN is expected to summarize the prefix with relevant features for the specific task. Such recurrent architectures are difficult to train due in particular to the problem of vanishing and exploding gradients [3]. This problem is partially solved with long

short-term memory (LSTM) units [4], which are crucial components in many state of the art architectures for tasks including handwriting recognition [5, 6], speech recognition [7, 8], image captioning [9] or machine translation [10].

We implemented and trained a LSTM RNN that reads the trajectory one GPS point at a time from the beginning to the end of each input prefix. This architecture is represented in Figure 2. Furthermore, in order to help the network to better identify short-term dependencies (such as the velocity of the taxi as the difference between two successive data points), we also considered a variant in which the input of the RNN is not anymore a single GPS point but a window of 5 successive GPS points of the prefix. The window shifts along the prefix by one point at each RNN time step.

3.2 Bidirectional Recurrent Neural Networks

We noticed that the most relevant parts of the prefix are its beginning and its end and we therefore tried a bidirectional RNN [11] (BRNN) to focus on these two particular parts. Our previously described RNN reads the prefix forwards from the first to the last known point, which leads to a final internal state containing more information about the last points (the information about the first points is more easily forgotten). In the BRNN architecture, one RNN reads the prefix forwards while a second RNN reads the prefix backwards. The two final internal states of the two RNNs are then concatenated and fed to a standard MLP that will predict the destination in the same way as in our previous models. The concatenation of these two final states is likely to capture more information about the beginning and the end of the prefix. Figure 3 represents the BRNN component of our architecture.

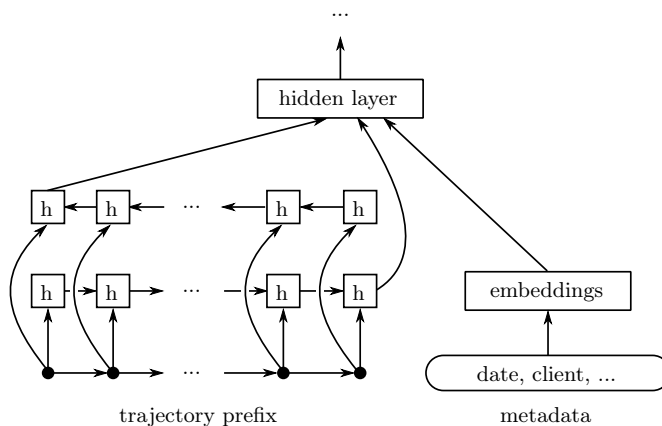


Fig. 3: Bidirectional RNN architecture.

3.3 Memory Networks

Memory networks [12] have been recently introduced as an architecture that can exploit an external database by retrieving and storing relevant information for each prediction. We have implemented a slightly related architecture, which is represented in Figure 4 and which we will now describe. For each prefix to predict, we extract m entire trajectories (which we call *candidates*) from the training dataset. We then use two neural network encoders (either feedforward or recurrent) to respectively encode the prefix and the candidates (same encoder for all the candidates). This results into $m + 1$ fixed-length representations in the same vector space so that they can be easily compared. Then we compute similarities by taking the dot products of the prefix representation with all the candidate representations. Finally we normalize these m similarity values with a softmax and use the resulting probabilities to weigh the destinations of the corresponding candidates. In other words, the final destination prediction of the prefix is the centroid of the candidate destinations weighted by the softmax probabilities. This is similar to the way we combine clusters in our previously described architectures.

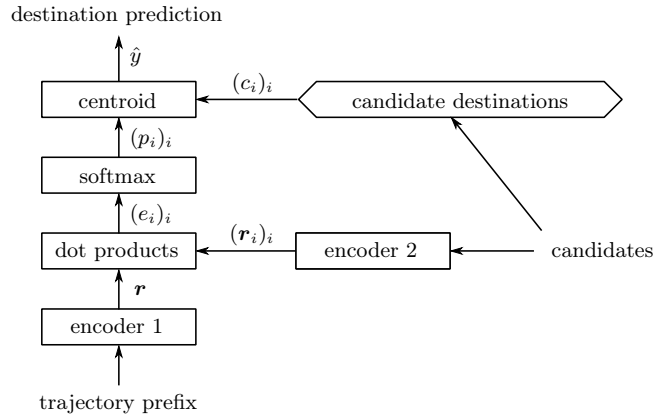


Fig. 4: Memory network architecture. The *encoder* are a generic bricks that take as input the trajectory points with metadata and process them through a feed-forward or recurrent neural net in order to return a fixed-size representations \mathbf{r} and $(\mathbf{r})_i$.

As the trajectory database is very large, such an architecture is quite challenging to implement efficiently. Therefore, for each prefix (more precisely for each batch of prefixes), we naively select $m = 10000$ random candidates. We believe that more sophisticated retrieving functions could significantly improve the results, but we did not have time to implement them. In particular, one could use a pre-defined (hand-engineered) similarity measure to retrieve the most similar candidates to the particular prefix.

The two encoders that map prefixes and candidates into the same representation space can either be feedforward or recurrent (bidirectional). As RNNs are more expensive to train (both in terms of computation time and RAM consumption), we had to limit ourselves to a MLP with one single hidden layer of 500 ReLUs for the encoder. We trained the architecture with a batch size of 5000 examples, and for each batch we randomly pick 10000 candidates from the training set.

4 Experimental Results

4.1 Custom Validation Set

As the competition testing dataset is particularly small, we can not reliably compare models on it. Therefore, for the purpose of this paper, we will compare our models on two bigger datasets: a validation dataset composed of 19427 trajectories and a testing dataset composed of 19770 trajectories. We obtained these new testing and validation sets by extracting (and removing) random portions of the original training set. The validation dataset is used to early-stop our training algorithms for each model based on the best validation score, while the testing dataset is used to compare our different trained models.

4.2 Results

Table 2 shows the testing scores of our various models on our custom testing dataset as well as on the competition ones. The different hyperparameters of each model have been tuned.

Table 2: Testing errors of our models. Contrary to model 1, model 2 predicts the destination directly without using the clusters. Model 3 only uses the prefix as input without any metadata. Model 4 only uses metadata as input without the prefix. While the BRNN of model 6 takes a single GPS point at each time step, the BRNN of model 7 takes five consecutive GPS points at each time step.

Model	Custom Test	Kaggle Public	Kaggle Private
1 MLP, clustering (winning model)	2.81	2.39	1.87⁸
2 MLP, direct output	2.97	3.44	3.88
3 MLP, clustering, no embeddings	2.93	2.64	2.17
4 MLP, clustering, embeddings only	4.29	3.52	3.76
5 RNN	3.14	2.49	2.39
6 Bidirectional RNN	3.01	2.65	2.33
7 Bidirectional RNN with window	2.60	3.15	2.06
8 Memory network	2.87	2.77	2.20
Second-place team	/	2.36	2.09
Third-place team	/	2.45	2.11
Average competition scores ⁹	/	3.02	3.11

Table 3 shows the p -values of Z -tests under the null hypotheses that our models perform as well two by two on our custom test dataset (more precisely the null hypothesis is that the error distributions are the same for each pair of models). The low p -values show that our scores are significant. This is not true for the competition test set, for which the test score differences are not always significant.

Table 3: p -values of Z -tests under the null hypotheses that the models have the same error distribution two by two on our custom test dataset.

Model	1	2	3	5	6	7
2	$4.8 \cdot 10^{-3}$					
3	$4.0 \cdot 10^{-2}$	$2.5 \cdot 10^{-1}$				
5	$1.4 \cdot 10^{-6}$	$5.1 \cdot 10^{-3}$	$1.6 \cdot 10^{-3}$			
6	$2.2 \cdot 10^{-3}$	$3.0 \cdot 10^{-1}$	$1.4 \cdot 10^{-1}$	$2.9 \cdot 10^{-2}$		
7	$3.8 \cdot 10^{-4}$	$1.4 \cdot 10^{-11}$	$6.5 \cdot 10^{-8}$	$0.0 \cdot 10^{+0}$	$4.2 \cdot 10^{-11}$	
8	$1.9 \cdot 10^{-1}$	$3.4 \cdot 10^{-2}$	$1.6 \cdot 10^{-1}$	$1.9 \cdot 10^{-5}$	$1.6 \cdot 10^{-2}$	$2.1 \cdot 10^{-6}$

5 Analysis of the results

Our winning ECML/PKDD challenge model is the MLP model that uses clusters of the destinations. However it is not our best model on our custom test, which is the BRNN with window. As the p -values in Table 3 suggest, the scores on our custom test set are significantly more confident than on the competition test set and we can therefore assert that our overall best model is the BRNN with window.

The results and p -values also prove that embeddings and clusters significantly improve our models. The importance of embeddings can also be confirmed by visualizing them. Figure 5 shows 2D t-SNE [13] projections for two of these embeddings and clear patterns can be observed, proving that quarters of hour and weeks of the year are important features for the prediction.

The reported score of the memory network is lower than the others but this might be due to the fact that it was not trained until convergence (we stopped after one week on a high end GPU).

The scores on our custom test set are higher than the scores on the public and private test set used for the competition. This suggests that the competition testing set is composed of rides that took place at very specific dates and times with very particular trajectory distributions. The gap between the public and private test sets is probably due to the fact that their size is particularly small. In contrast, our validation and test sets are big enough to obtain more significant statistics.

⁸ our winning submission on Kaggle scored 2.03 but the model had not been trained until convergence

⁹ average over 381 teams, the submissions with worse scores than the public benchmark (in which the center of Porto is always predicted) have been discarded

All the models we have explored are very computationally intensive and we thus had to train them on GPUs to avoid weeks of training. Our competition winning model is the least intensive and can be trained in half a day on GPU. On the other hand, our recurrent and memory networks are much slower and we believe that we could reach even better scores by training them longer.

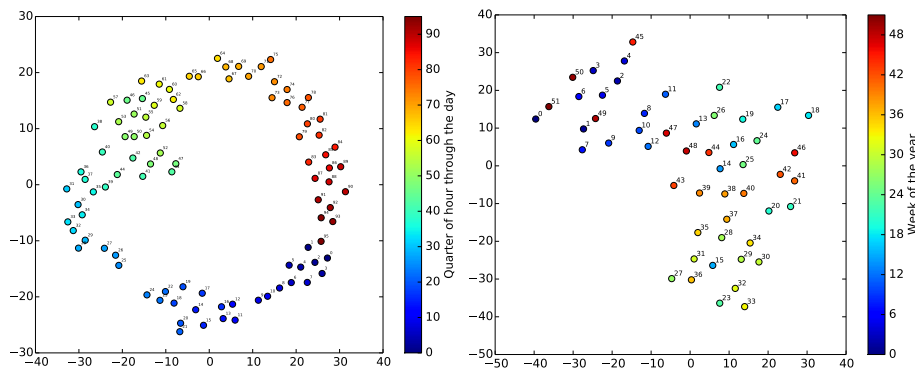


Fig. 5: t-SNE 2D projection of the embeddings. **Left:** quarter of hour at which the taxi departed (there are 96 quarters in a day, so 96 points, each one representing a particular quarter). **Right:** week of the year during which the taxi departed (there are 52 weeks in a year, so 52 points, each one representing a particular week).

Conclusion

We introduced an almost fully-automated neural network approach to predict the destination of a taxi based on the beginning of its trajectory and associated metadata. Our best model uses a recurrent bidirectional neural network to encode the prefix, several embeddings to encode the metadata and destination clusters to generate the output.

One potential limitation of our clustering-based output layer is that the final prediction can only fall in the convex hull of the clusters. A potential solution would be to learn the clusters as parameters of the network and initialize them either randomly or from the mean-shift clusters.

Concerning the memory network, one could consider more sophisticated ways to extract candidates, such as using an hand-engineered similarity measure or even the similarity measure learnt by the memory network. In this latter case, the learnt similarity should be used to extract only a proportion of the candidates in order let a chance to candidates with poor similarities to be selected. Furthermore, instead of using the dot product to compare prefix and candidate representations, more complex functions could be used (such as the concatenation of the representations followed by non-linear layers).

References

1. Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 315–323, 2011.
2. Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *The Journal of Machine Learning Research*, 3:1137–1155, 2003.
3. Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *Neural Networks, IEEE Transactions on*, 5(2):157–166, 1994.
4. Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
5. Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and Jürgen Schmidhuber. A novel connectionist system for unconstrained handwriting recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 31(5):855–868, 2009.
6. Patrick Doetsch, Michal Kozielski, and Hermann Ney. Fast and robust training of recurrent neural networks for offline handwriting recognition. In *Frontiers in Handwriting Recognition (ICFHR), 2014 14th International Conference on*, pages 279–284. IEEE, 2014.
7. Alan Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6645–6649. IEEE, 2013.
8. Haşim Sak, Andrew Senior, and Françoise Beaufays. Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition. *arXiv preprint arXiv:1402.1128*, 2014.
9. Kelvin Xu, Jimmy Ba, Ryan Kiros, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. *arXiv preprint arXiv:1502.03044*, 2015.
10. Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
11. Alex Graves, Santiago Fernández, and Jürgen Schmidhuber. Bidirectional lstm networks for improved phoneme classification and recognition. In *Artificial Neural Networks: Formal Models and Their Applications–ICANN 2005*, pages 799–804. Springer, 2005.
12. Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. *arXiv preprint arXiv:1410.3916*, 2014.
13. Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.