

Abstract

We summarise previous work showing that the basic sigmoid activation function arises as an instance of Bayes's theorem, and that recurrence follows from the prior. We derive a layer-wise recurrence without the assumptions of previous work, and show that it leads to a standard recurrence with modest modifications to reflect use of log-probabilities. The resulting architecture closely resembles the Li-GRU which is the current state of the art for ASR. Although the contribution is mainly theoretical, we show that it is able to outperform the state of the art on the TIMIT and AMI datasets.

Context and motivation

In signal processing, modelling of context can be achieved via recurrence. The most successful architectures are based on the long short-term memory (LSTM), which uses a memory cell and different gates to filter out irrelevant information. Recent efforts have been pursued to reduce the size of recurrent units and avoid redundancies. This notably lead to the gated recurrent unit (GRU) [1] in 2014 and light GRU [3] in 2018. Following the recent work of Garner and Tong [2], our motivation is to provide a probabilistic interpretation of recurrent units and shed some light on the seemingly ad-hoc concepts of gates and memory cells.

Bayesian interpretation of recurrence

Consider an input sequence $\mathbf{X}_T = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{F \times T}$ of length T , where each observation \mathbf{x}_t is a vector of size F . We assume that there are H hidden features $\{\phi_i | i = 1, \dots, H\}$ that we wish to detect along the sequence. At each timestep t , a feature has two possible states: present or absent, that we write as $\phi_{t,i}$ and $\neg\phi_{t,i}$ respectively. We want to build a layer of H recurrent units that will output the stacked probabilities $\mathbf{h}_t := P(\phi_t | \mathbf{X}_t) \in [0, 1]^H$ of the different features being present at each timestep $t = 1, \dots, T$.

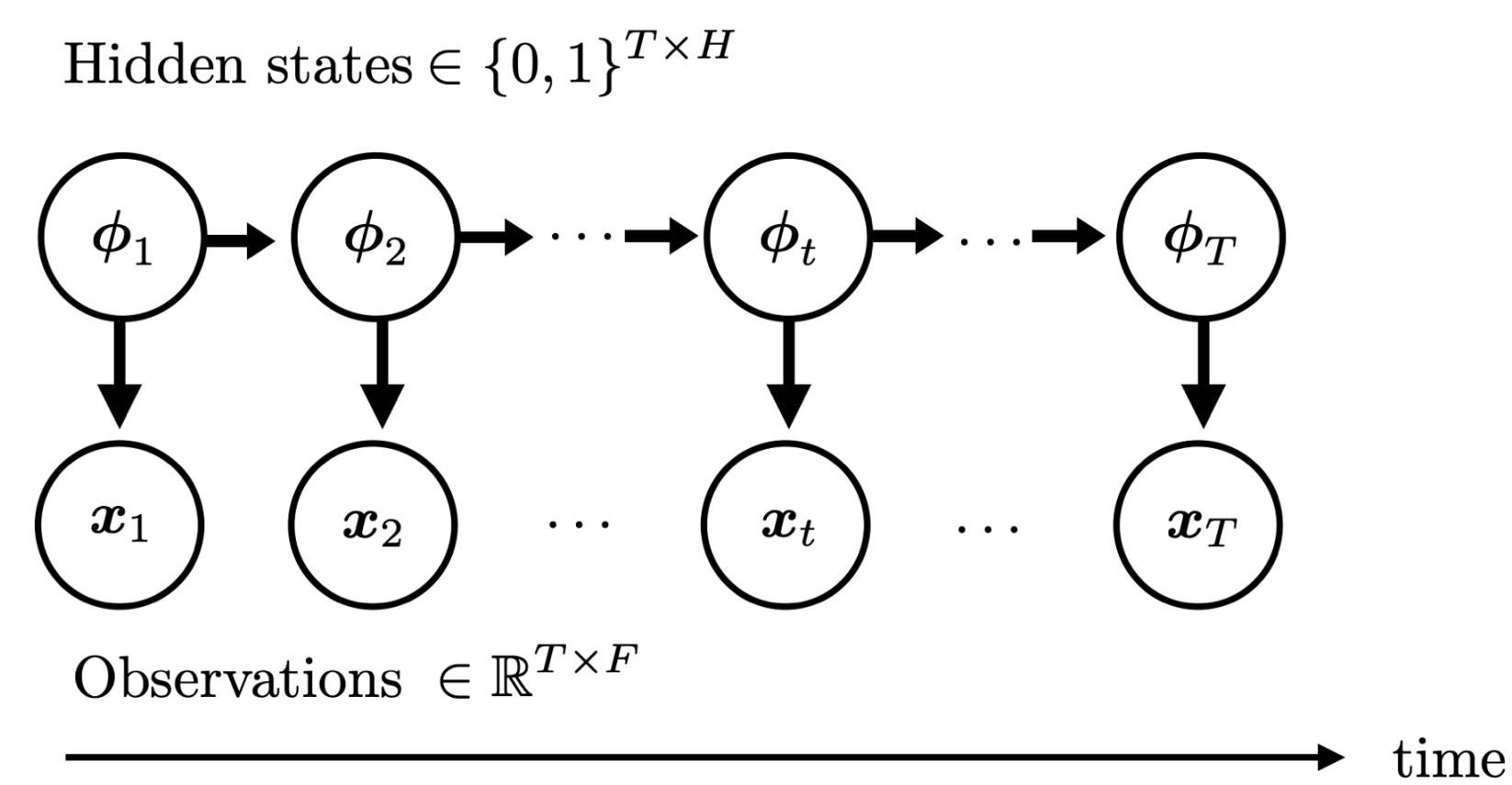


Figure 1. Diagram of the approach

Bayes's theorem

For the two-class case, Bayes's formula can be rewritten in a sigmoid form for the whole layer,

$$\mathbf{h}_t = \sigma \left[\log(\mathbf{r}_t) + \text{logit}(\mathbf{p}_t) \right]. \quad (1)$$

The posterior \mathbf{h}_t , representing the desired unit outputs, is expressed as a function of the ratio of likelihood \mathbf{r}_t and prior \mathbf{p}_t , defined as

$$\mathbf{r}_t := \frac{p(\mathbf{x}_t | \phi_t)}{p(\mathbf{x}_t | \neg\phi_t)} \quad \text{and} \quad \mathbf{p}_t := P(\phi_t | \mathbf{X}_{t-1}). \quad (2)$$

Ratio of likelihood and Prior

If we assume that the likelihood of observing \mathbf{x}_t given the current state of the features ϕ_t can be represented with multivariate normal distributions, the ratio of likelihood can be expressed as

$$\mathbf{r}_t = \exp \left[\mathbf{W}_r^T \mathbf{x}_t + \mathbf{b}_r \right]. \quad (3)$$

Assuming that the H feature probabilities $h_{t-1,i}$ of the layer are independently beta distributed, they can be combined into a single prior probability,

$$\mathbf{p}_t = \sigma \left[\mathbf{V}_p \log(\mathbf{h}_{t-1}) + \mathbf{b}_p \right]. \quad (4)$$

Resulting forward pass

By plugging equations (3) and (4) into (1), we get the following update equation,

$$\mathbf{h}_t = \sigma \left[\mathbf{W}_h \mathbf{x}_t + \mathbf{V}_h \log(\mathbf{h}_{t-1}) + \mathbf{b}_h \right], \quad (5)$$

or in log-probabilities form, i.e. $\mathbf{h}_t := \log [P(\phi_t | \mathbf{X}_t)]$,

$$\mathbf{h}_t = \text{softplus} \left[\mathbf{W}_h \mathbf{x}_t + \mathbf{V}_h \mathbf{h}_{t-1} + \mathbf{b}_h \right], \quad (6)$$

where \mathbf{W}_h , \mathbf{V}_h and \mathbf{b}_h are representative of the distributions of \mathbf{x}_t and \mathbf{h}_{t-1} when the features are present or absent, and can be treated as trainable parameters of the model.

Adding a context relevance gate

We define a binary state variable $\rho_{t,i}$, which indicates if \mathbf{x}_t is relevant for the occurrence of $\phi_{t,i}$. The associated probabilities $z_{t,i} = P(\rho_{t,i} | \mathbf{X}_t)$ can be computed with equation (5). The desired output probabilities $h_{t,i}$ can then be expressed by marginalizing over $\rho_{t,i}$ so that we get the forward pass of a Li-GRU with minor modifications,

$$\mathbf{z}_t = \sigma \left[\mathbf{W}_z \mathbf{x}_t + \mathbf{V}_z \log(\mathbf{h}_{t-1}) + \mathbf{b}_z \right] \quad (7a)$$

$$\tilde{\mathbf{h}}_t = \sigma \left[\mathbf{W}_h \mathbf{x}_t + \mathbf{V}_h \log(\mathbf{h}_{t-1}) + \mathbf{b}_h \right] \quad (7b)$$

$$\mathbf{h}_t = \mathbf{z}_t * \tilde{\mathbf{h}}_t + (1 - \mathbf{z}_t) * \mathbf{h}_{t-1}. \quad (7c)$$

Additionally, the input of the i -th layer corresponds to the log-probabilities of the previous layer $\mathbf{x}_t^{[i]} = \log(\mathbf{h}_t^{[i-1]})$.

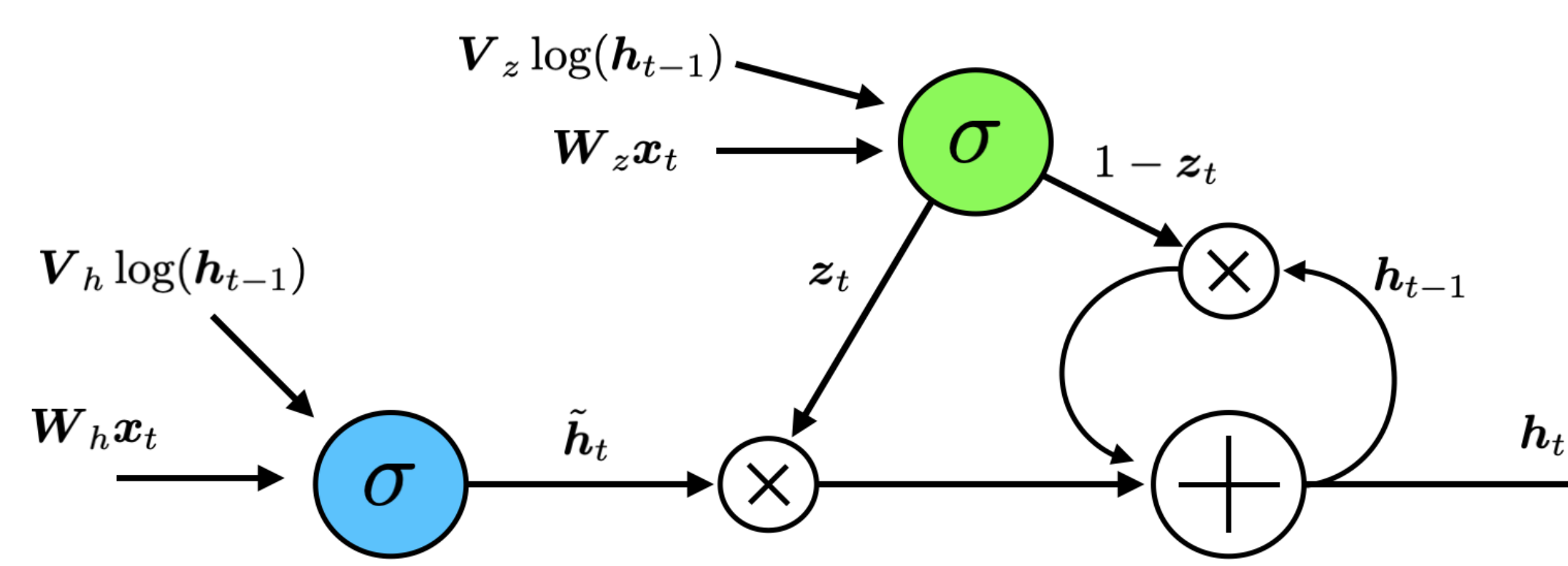


Figure 2. Diagram of the Li-BRU defined by equations (7)

Comparison with Li-GRU

The forward pass of a Li-GRU is the following,

$$\mathbf{z}_t = \sigma \left[\mathbf{W}_z \mathbf{x}_t + \mathbf{V}_z \mathbf{h}_{t-1} + \mathbf{b}_z \right] \quad (8a)$$

$$\tilde{\mathbf{h}}_t = \text{ReLU} \left[\mathbf{W}_h \mathbf{x}_t + \mathbf{V}_h \mathbf{h}_{t-1} + \mathbf{b}_h \right] \quad (8b)$$

$$\mathbf{h}_t = \mathbf{z}_t * \tilde{\mathbf{h}}_t + (1 - \mathbf{z}_t) * \mathbf{h}_{t-1}. \quad (8c)$$

The differences with our unit are

- the ReLU activation, which is an approximation of the softplus
- the gate applied to log-probabilities.

References

- K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 EMNLP Conference*, page 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics.
- P. N. Garner and S. Tong. A Bayesian approach to recurrence in neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- M. Ravanelli, A. Bordes, and Y. Bengio. Light gated recurrent units for speech recognition. *Transactions on Emerging Topics in Computational Intelligence*, 2(2):92–102, 2018.

Experiments

We perform a self-consistent comparison between recurrent units on 2 speech recognition tasks, using the TIMIT and AMI corpora. Following the pytorch-kaldi implementation of [3], all presented experiments use a recurrent architecture with 4 layers of H=550 bidirectional units. The F=50 fMLLR input features are extracted via the Kaldi recipe.

Without the update gate

We first test the simple BRU from equation (6) on TIMIT and compare it to conventional units of the same size. Units that have a feedback on log-probabilities perform significantly better.

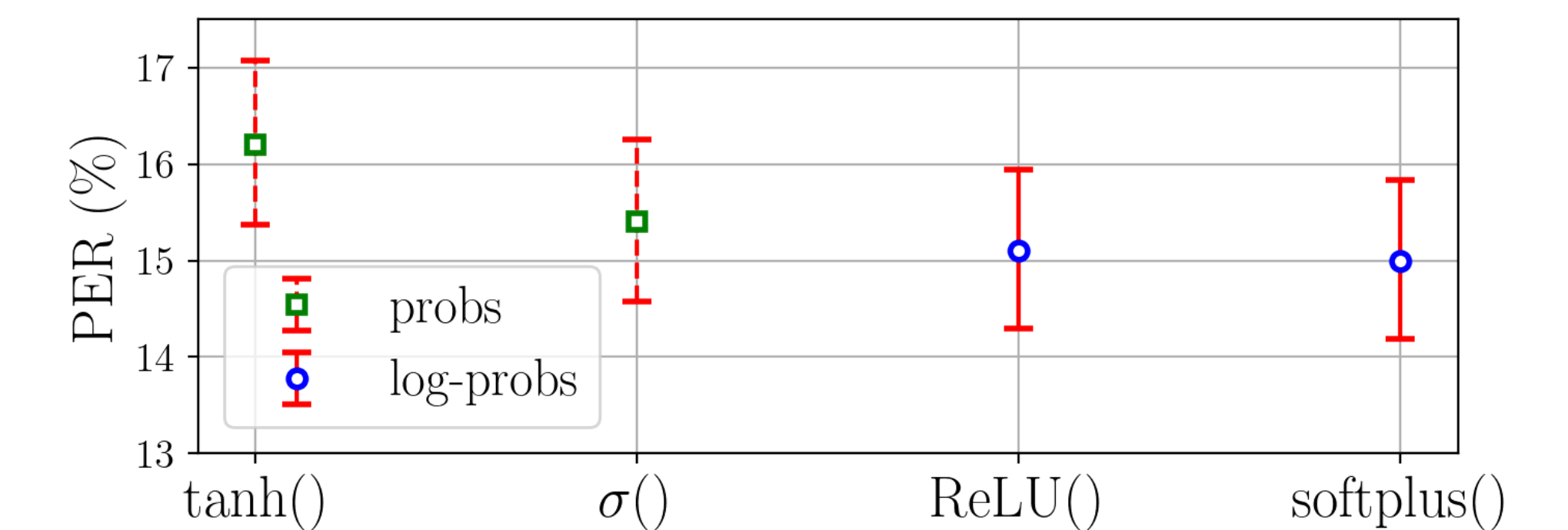


Figure 3. Error-rates on TIMIT testset for various RNN architectures.

With the update gate

We now test the complete Li-BRU from equations (7) to state of the art recurrent units.

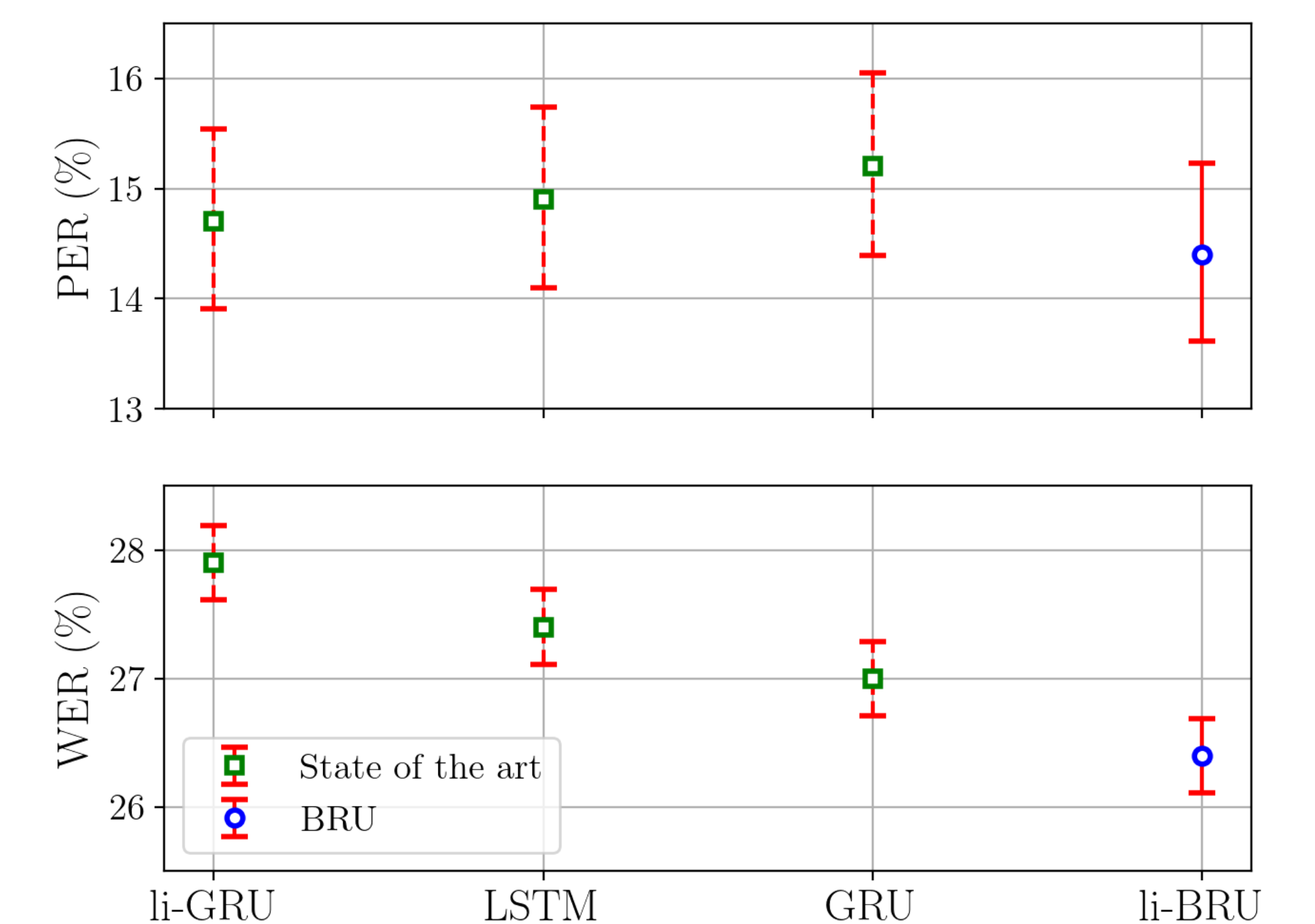


Figure 4. Error-rates on TIMIT (top) and AMI (bottom) for various RNN architectures.

The Li-BRU outperforms all other state of the art recurrent architectures on both TIMIT and AMI datasets with error-rates of 14.4 %, and 26.4% respectively. We found that applying the gate to probabilities or log-probabilities is practically equivalent. The root of improvement of the Li-BRU over the Li-GRU therefore seems to lie in the use of the softplus function instead of its approximation by the ReLU.

Conclusion

In previous work, it was shown that a Bayesian analysis of a sigmoid activation function led naturally to a unit-wise recurrence and, with approximations, to a layer-wise recurrence. In this paper, in a mainly theoretical contribution, we have shown that beta-distributed sigmoid outputs feeding into another sigmoid unit constitute a layer-wise recurrence without approximations. Without a forget gate, this reduces to a standard fully-connected recurrence, but with a softplus activation. Given that the update gate of a GRU can also be derived probabilistically, the approach led naturally to comparison with a Li-GRU. In an experimental evaluation, we confirmed that the resulting light Bayesian recurrent unit (Li-BRU) can modestly but significantly outperform the state of the art on two ASR tasks (TIMIT and AMI datasets), demonstrating the importance of the probabilistic derivation. More generally, the new techniques contribute to a growing toolkit of Bayesian approaches for neural architectures.