

Structured Sparsification of Gated Recurrent Neural Networks

Ekaterina Lobacheva^{*1} Nadezhda Chirkova^{*1} Alexander Markovich² Dmitry Vetrov^{1,3}

¹ Samsung-HSE Laboratory, Higher School of Economics;

² Higher School of Economics; ³ Samsung AI Center Moscow



SAMSUNG
Research



HIGHER SCHOOL OF ECONOMICS
NATIONAL RESEARCH UNIVERSITY

Moscow, Russia

Overview

Goal: learn structured sparsity in gated RNNs (LSTM, GRU, ...)

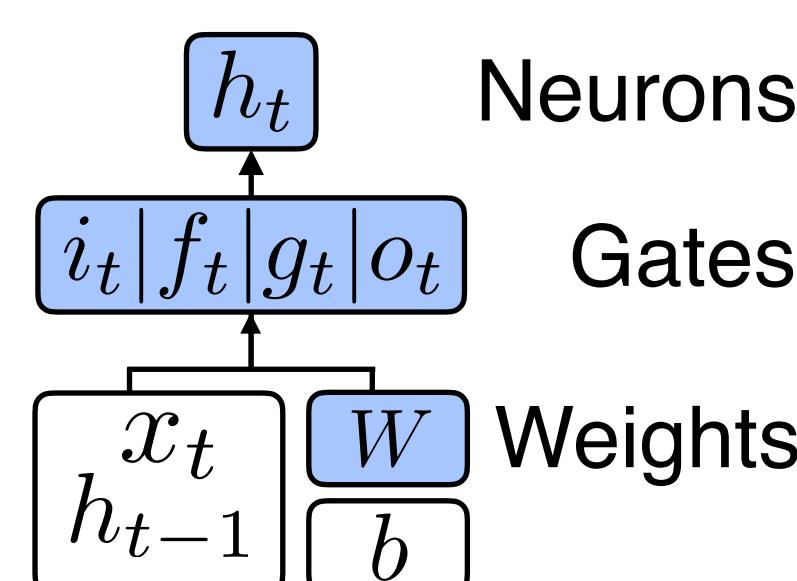
Structured sparsity in RNNs:

1. Learn sparse weight matrices [1] → reduce model size
2. Remove neurons by removing all the weights associated with these neurons [2] → faster prediction

Our approach:

add intermediate level of sparsification — → simplify LSTM cell & reveal specifics of LSTM workflow

3 sparsity levels:



Results:

- Improved neuron-wise compression
- LSTM uses different gate structures to solve different tasks and to model long / short dependencies

Main idea

Remove weights by groups corresponding to gates / neurons:

$$\begin{bmatrix} i_t \\ f_t \\ g_t \\ o_t \end{bmatrix} = \begin{array}{l} \text{sigm} \\ \text{sigm} \\ \tanh \\ \text{sigm} \end{array} \begin{bmatrix} W^x & W^h & b \end{bmatrix} \begin{bmatrix} x_t \\ h_{t-1} \end{bmatrix} + \begin{bmatrix} h_t \\ o_t \end{bmatrix}$$

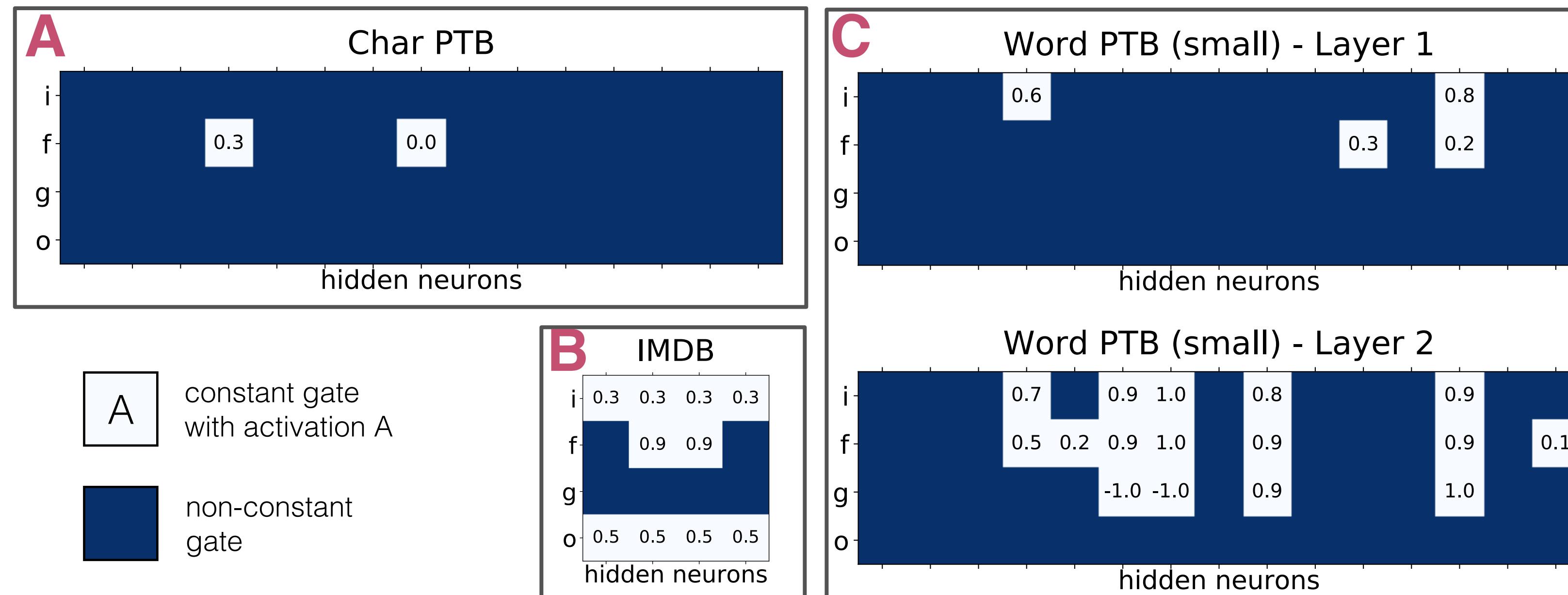
zero column
no outgoing edges
↓
remove neuron

zero row
no ingoing edges
↓
constant gate

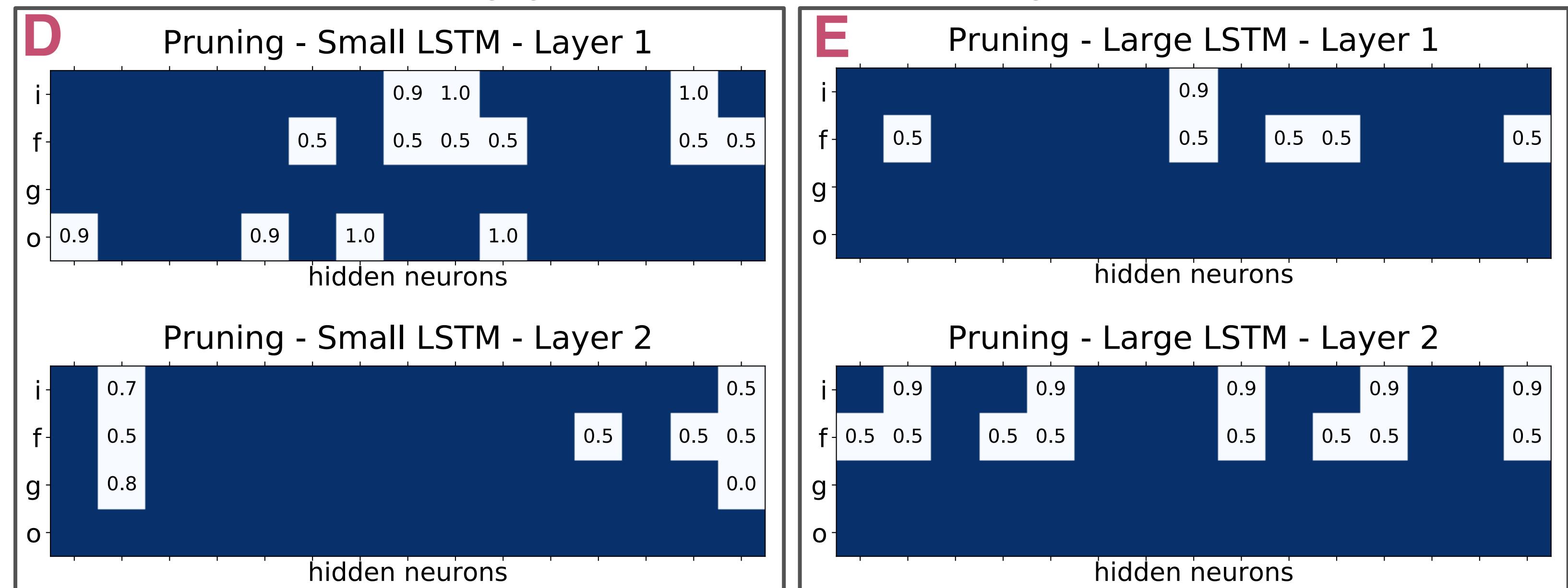
$c_t = f_t \odot c_{t-1} + i_t \odot g_t$

Experiments: qualitative

Resulting gate structures: Bayes W+G+N



Resulting gate structures: Pruning W+G+N



Resulting gate structures vary for different tasks: compare A, B, C

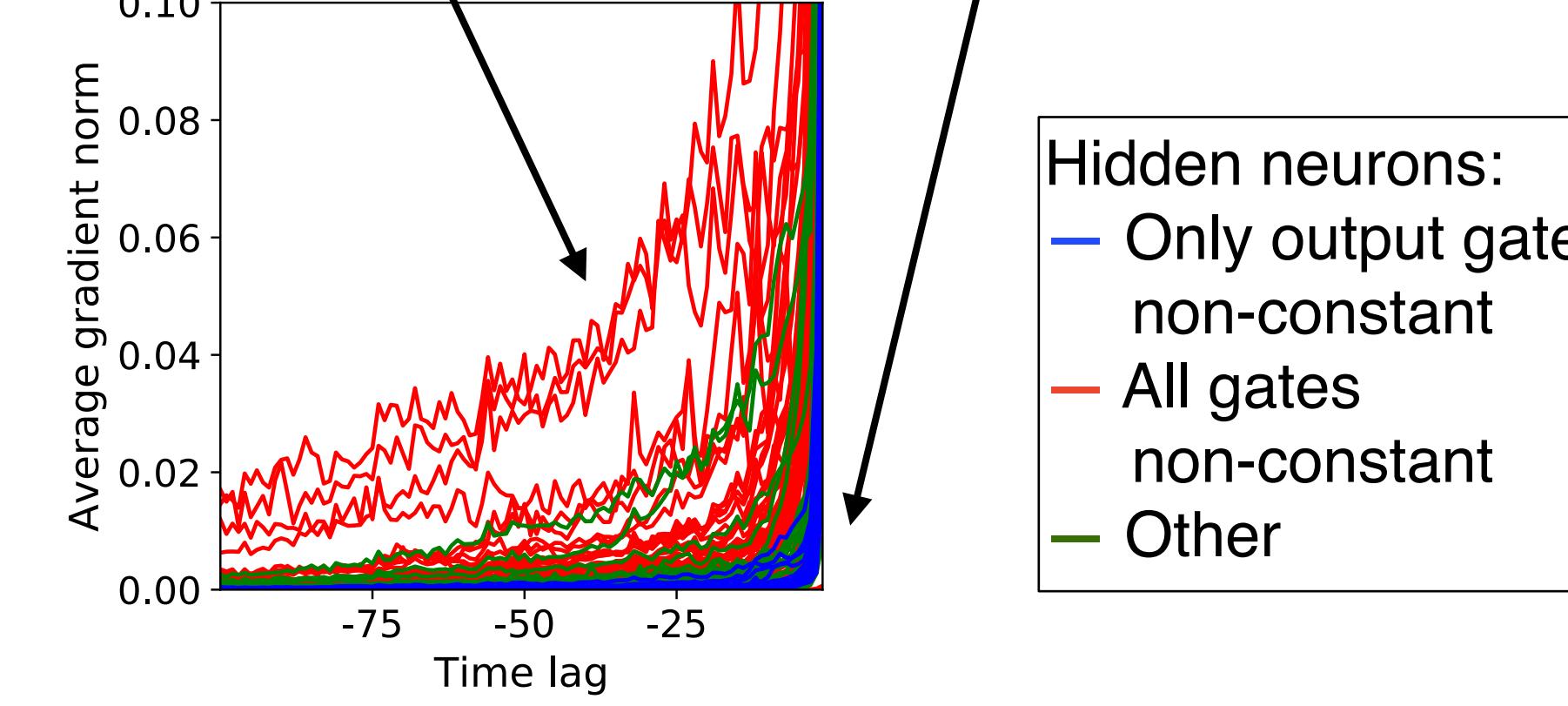
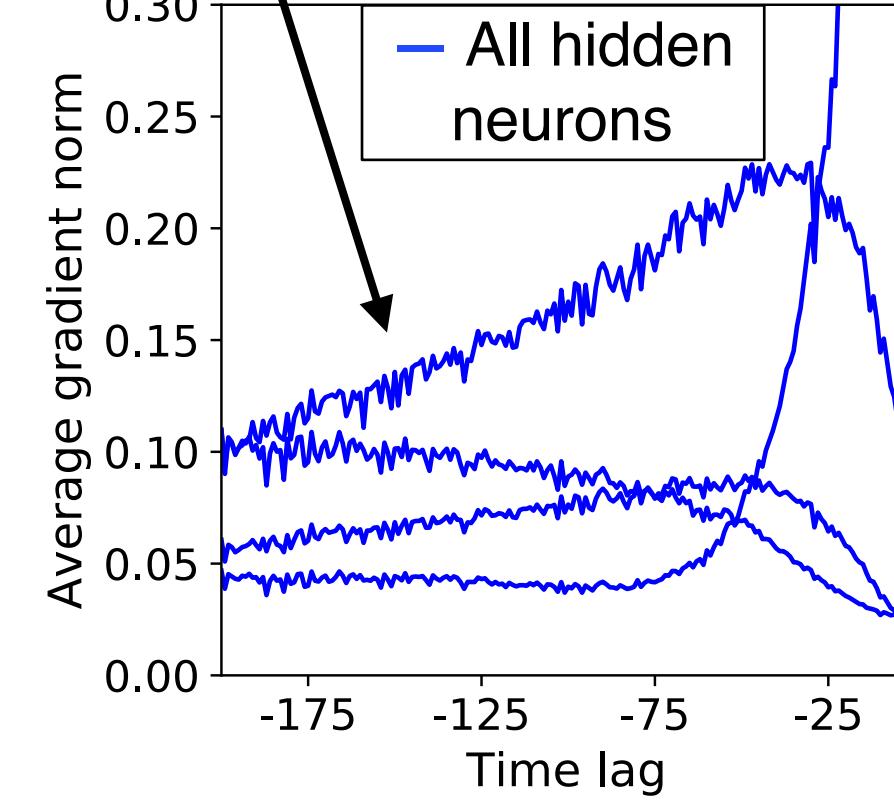
char-level language model (A): all gates non-constant

classification (B): non-constant g — long dependencies are highly important

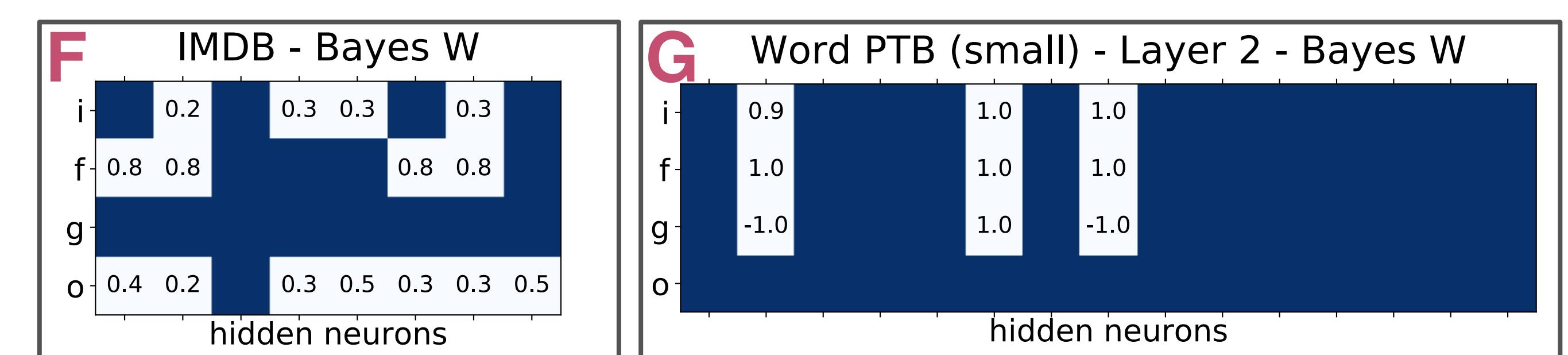
word-level language model, 2nd layer (C): all gates non-constant — full LSTM units — long dependencies

only non-constant o — vanilla recurrent units — short dependencies

Averaged over validation sequences, norm of the gradients of hidden neurons w. r. t. LSTM input for different time-lag



Resulting gate structures: Bayes W



Gate structure intrinsically exists in LSTM:

- for the same task, similar gate structures are obtained in different sparsification frameworks (compare C and D)
- characteristic patterns in gate structures are obtained even with unstructured sparsification (compare B and F; C and G), but with lower overall compression

References

[1] Chirkova, N. et al.

Bayesian compression for natural language processing. EMNLP 2018

[2] Wen, W. et al.

Learning intrinsic sparse structures within long short-term memory. ICLR 2018

[3] Louizos, C. et al.

Bayesian compression for deep learning. NeurIPS 2017

[4] Molchanov D. et al.

Variational dropout sparsifies deep neural networks. ICML 2017