

# A Pruned RNNLM Lattice-Rescoring Algorithm for Automatic Speech Recognition

Hainan Xu<sup>1,2</sup>, Tongfei Chen<sup>1</sup>, Dongji Gao<sup>1</sup>, Yiming Wang<sup>1</sup>, Ke Li<sup>1</sup>, Nagendra Goel<sup>3</sup>, Yishay Carmiel<sup>2</sup>, Daniel Povey<sup>1</sup>, Sanjeev Khudanpur<sup>1</sup>

<sup>1</sup>Center for Language and Speech Processing, Johns Hopkins University; <sup>2</sup>IntelligentWire, Seattle WA; <sup>3</sup>Go-Vivace Inc., USA

## Overview

- Usually lattice-rescoring uses  $n$ -gram approximation to limit search space;
- We propose a heuristic that finds more promising arcs to expand, and use it for pruning;
- Complexity of the algorithm grows approximately (empirically) linear with  $n$ -gram order, compared with exponential growth of the baseline algorithm;
- 4X and 10X faster for 4-gram and 5-gram;
- The heuristic also consistently improves WER;
- The evaluation is done with TensorFlow RNNLMs. We open source the integration of TensorFlow to Kaldi.

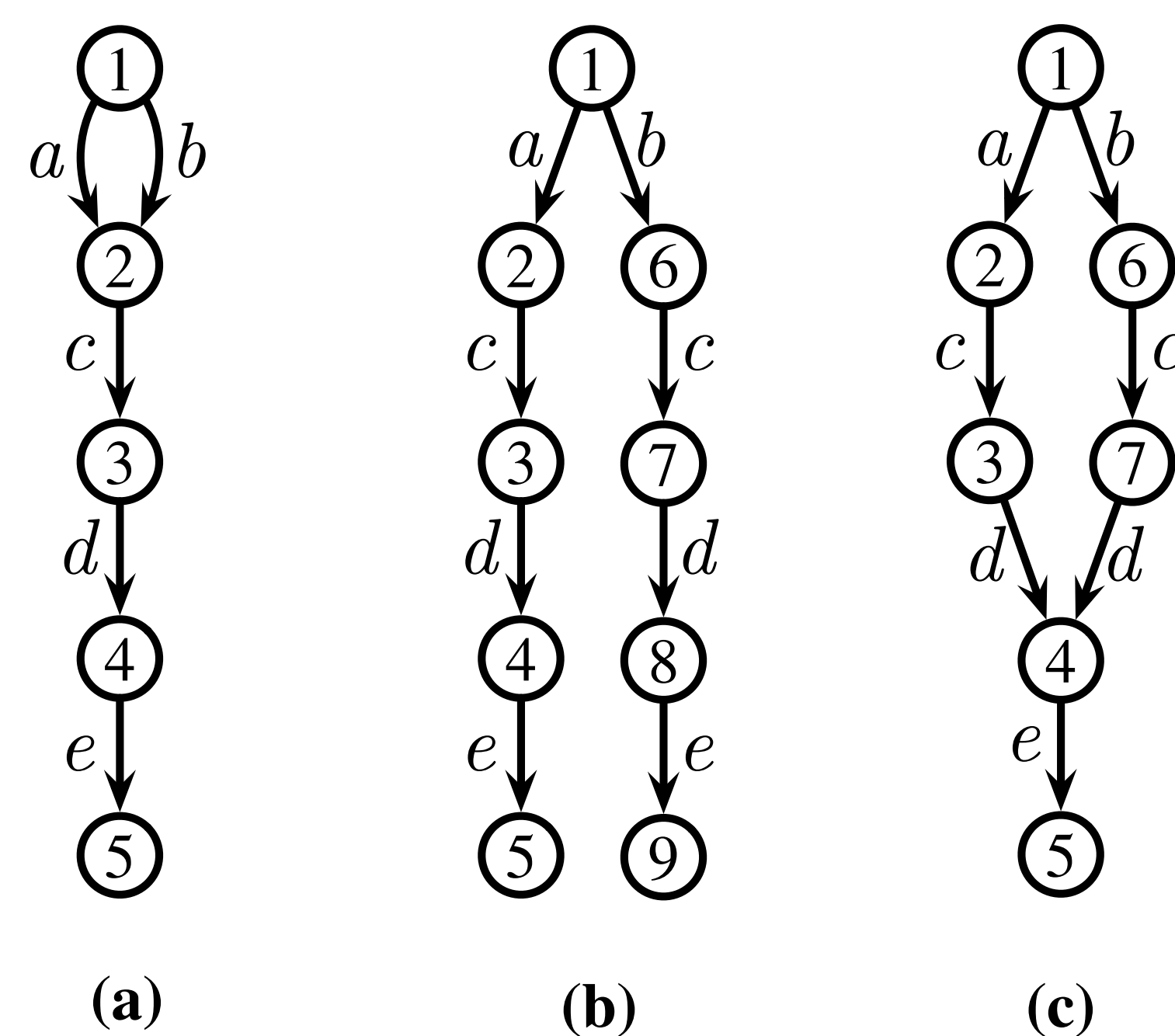
## Lattice Rescoring

- In speech recognition, decoding is usually done on a static decoding graph compiled from an  $n$ -gram;
- RNNLM rescoring helps further reduce WERs by (partially) replacing LMs weights on a decoded lattice;
- A naive implementation to rescore the lattice is too costly – it runs exponentially w.r.t. lattice-depth;
- An  $n$ -gram approximation algorithm is commonly used in order to limit the search space.

## Contact Information

- Web: <http://www.hainanxv.com>
- Email: [hxu31@jhu.edu](mailto:hxu31@jhu.edu)

## Analysis of Old Algorithm



- In 3-gram approximation, states 4 and 8 in (b) are merged as state 4 in (c);
- state 4 encodes history of either (a, c, d) or (b, c, d). The choice is arbitrary, and affects the weight computed for  $p(e | 4)$ .

## Pruned Algorithm

- For each arc to be expanded, we compute a score reflecting how likely this arc will become part of the best-path;
- Arcs that are not very promising (out of the beam) are not expanded;
- Arcs that are more promising get expanded first, so that output lattice states encode “better” history.

## Heuristic

- The heuristic is computed as

$$H(c) = \alpha(c) + \beta(a) + \delta(c) \quad (1)$$

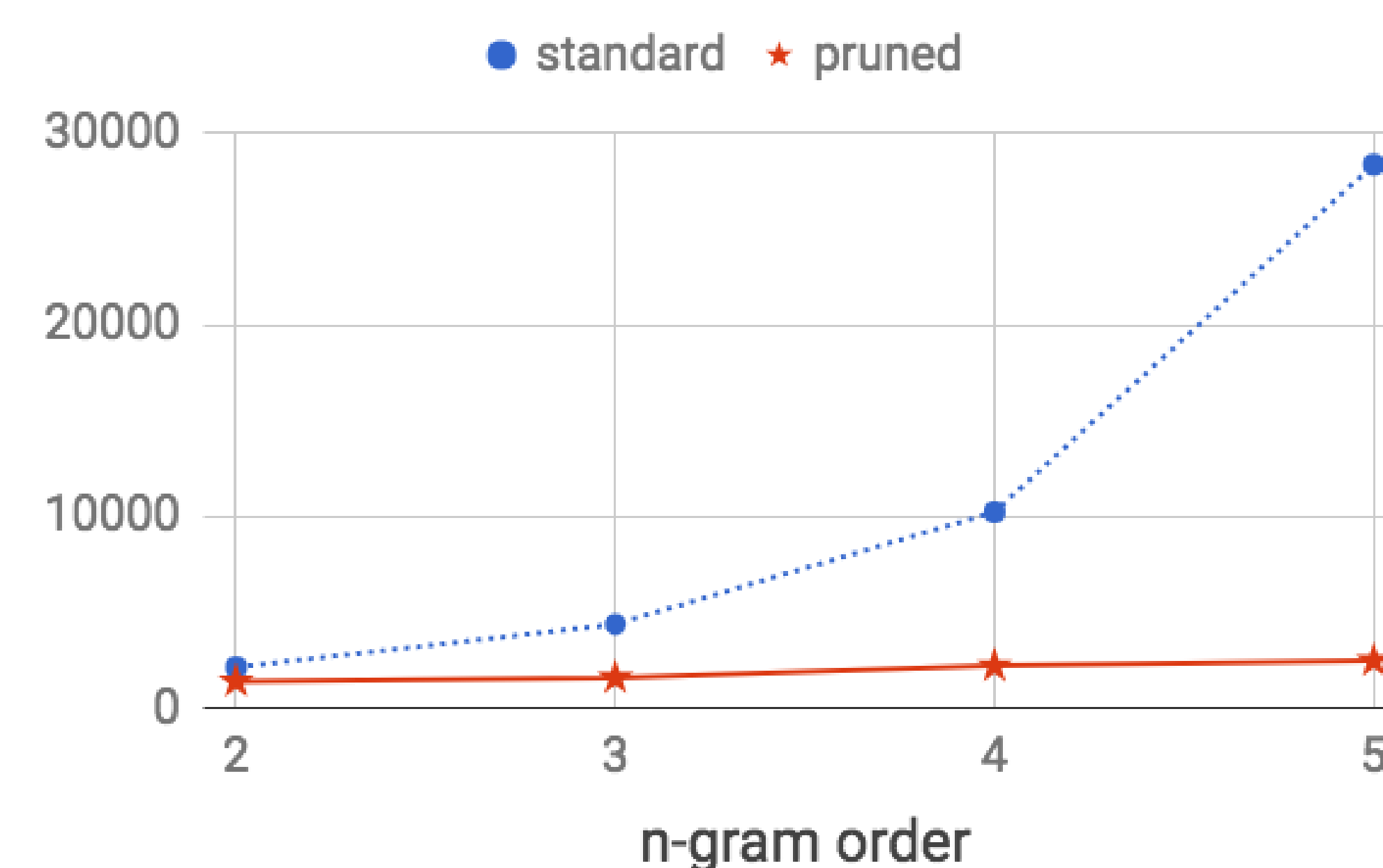
- $c$ : a state in the output lattice;
- $a$ : the corresponding state in the input lattice;

- $\alpha(c)$  is the forward-cost for  $c$  in the output lattice
- $\beta(a)$  is the backward-cost for  $a$  in the input lattice
- $\delta(c)$  is an “expectation” of  $\beta(c) - \beta(a)$

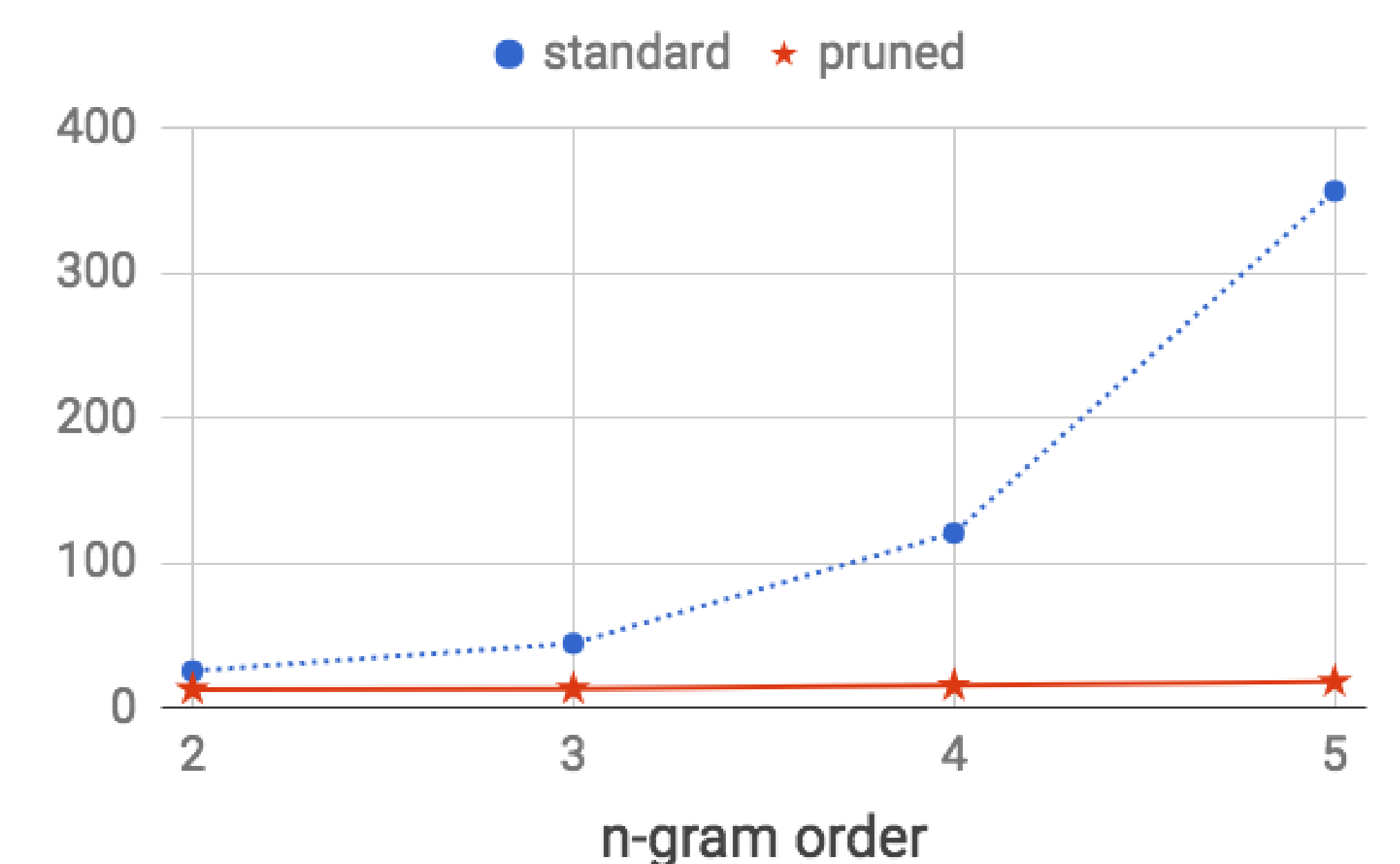
$$\delta(c) = \begin{cases} \beta(c) - \beta(a), & \beta(c) < +\infty \\ \delta(\text{prev}(c)), & \beta(c) = +\infty \end{cases} \quad (2)$$

$\text{prev}(c)$  is the previous state of  $c$  on the best path from start to  $c$ .

## Lattice-rescoring Speed



## Output Lattice Size (arcs per frame)



## Acknowledgements

This work was partially supported by DARPA LORELEI award number HR0011-15-2-0024, NSF Grant No CRI-1513128 and IARPA MATERIAL award number FA8650-17-C-9115 and by IntelligentWire. The authors would also like to thank the TensorFlow team at Google for their help during the project.

## Word-error-rate

Corpus	Test set	ARPA baseline	RNNLM rescoring with $n$ -gram approximation					
			2-gram		3-gram		4-gram	
			standard	pruned	standard	pruned	standard	pruned
AMI-IHM (0.5)	dev	24.2	24.5	<b>24.0</b>	23.7	<b>23.4</b>	23.4	<b>23.3</b>
	eval	25.4	25.8	<b>25.0</b>	24.6	<b>24.4</b>	24.3	<b>24.2</b>
SWBD (0.8)	swbd	8.1	8.6	<b>8.2</b>	7.4	<b>7.2</b>	7.2	<b>7.1</b>
	eval2000	12.4	12.9	<b>12.3</b>	11.7	<b>11.5</b>	11.5	<b>11.3</b>
WSJ (0.8)	dev93	7.6	7.2	<b>6.9</b>	6.4	<b>6.2</b>	6.4	<b>6.2</b>
	eval92	5.1	4.6	<b>4.2</b>	4.1	<b>3.9</b>	3.9	<b>3.8</b>
LIB (0.5)	test-clean	6.0	5.5	<b>5.1</b>	4.9	<b>4.8</b>	4.8	<b>4.7</b>
	test-other	15.0	14.0	<b>13.2</b>	12.7	<b>12.4</b>	12.4	<b>12.3</b>
	dev-clean	5.7	5.0	<b>4.8</b>	4.4	<b>4.3</b>	4.3	<b>4.3</b>
	dev-other	14.5	13.7	<b>12.9</b>	12.3	<b>12.0</b>	11.9	<b>11.7</b>

Table 1: WER of Lattice-rescoring of Different RNNLMs