



# A Parallelizable Lattice Rescoring Strategy with Neural Language Models

Ke Li<sup>1</sup>, Daniel Povey<sup>2</sup>, Sanjeev Khudanpur<sup>1</sup>

<sup>1</sup>The Johns Hopkins University, USA, <sup>2</sup>Xiaomi Corp., China.

## Motivation

- ▶ Lattice rescoring usually involves lattice expansion
  - ▶ The general goal of lattice expansion is to make arcs on highly likely paths have unique histories
  - ▶  $n$ -gram approximation, which merges histories that share  $(n - 1)$  most recent words, may sacrifice accuracy and waste computation on less likely paths. **Whether can we do better?**
- ▶ A major speedup bottleneck of lattice rescoring with neural LMs is LM evaluation
  - ▶ Existing methods such as caching computed LM scores and pruning-based algorithms speedup the process by reducing the number of LM evaluations
  - ▶ While the sequential LM evaluation order in a lattice is still inefficient. **How to parallel LM evaluations within a lattice?**

## A Parallel Lattice Rescoring Strategy (“Non-iterative”)

### Rescoring Procedure

- ▶ Step 1 - Lattice Expansion: expand a lattice with a posterior-based method (with beam pruning applied beforehand)
- ▶ Step 2 - Lattice-to-List Conversion: convert the expanded lattice into a minimal list of hypotheses that cover every arc
- ▶ Step 3 - Score computation and estimation: compute LM scores of the lists in parallel and estimate LM scores for each arc
- ▶ Step 4 - Integrate scores back to the expanded lattice

### Posterior-based Lattice Expansion Algorithm

- ▶ Idea - only expand very possible arcs, e.g. arc posteriors  $> \epsilon$
- ▶ The basic question is whether an incoming arc should be split off from the rest of incoming arcs to its destination-state
- ▶ The rule is to allocate a new copy of the destination-state if the arc posterior  $> \epsilon$  (a predefined threshold), otherwise transition to the original destination-state

### Lattice-to-List Conversion

- ▶ Definition of *path cover*: a set of paths such that every arc in the lattice is covered by at least one path
- ▶ Lattice-to-List conversion aims to find a minimal path cover with condition that each path is the best one for at least one arc it has

### Estimation of Neural LM Scores

- ▶ Estimation of neural LM scores for each arc is needed since an arc may be shared by multiple paths
- ▶ We experiment with three approximation methods:
  - ▶ Simply average neural LM scores from shared paths
  - ▶ Perform weighted average with weights as neural LM scores of histories on shared paths
  - ▶ Choose the neural LM score from the lowest-cost path among the shared paths (Referred to “Semi-Viterbi” in experiments)

## Lattice-to-List Conversion

### Algorithm 1 A Constrained Path Cover Algorithm

**Input:**  $L$ : a lattice  
**Output:**  $O$ : a list of paths, each is represented as a linear FST.

```

1: procedure CONSTRAINEDPATHCOVER( $L$ )
2:   TopologicalSort( $L$ )
3:    $P \leftarrow []$            ▷ A list of pairs of a path and its cost
4:    $\alpha, \beta \leftarrow$  ViterbiForwardBackward( $L$ )
5:   for  $s = 0 : S - 1$  do           ▷ Loop over states
6:     for  $e \in s.out$  do           ▷ Loop over outgoing arcs of  $s$ 
7:       if best path including  $e$  is not generated then
8:          $p, c \leftarrow$  BestPathForAnArc( $\alpha, \beta, s, e$ )
9:          $P.append((p, c))$ 
10:  Sort( $P$ )           ▷ Sort paths based on their costs
11:   $O \leftarrow$  ConstructOutputLattice( $P$ )

```

## A Refined Parallel Lattice Rescoring Strategy (“Iterative”)

- ▶ Apply score replacement on top of the introduced non-iterative parallel lattice rescoring strategy
  - ▶ Score replacement means replacing  $n$ -gram scores with neural LM ones for lattices from first-pass decoding
  - ▶ It is referred to “iterative” since rescoring happens twice

## Experimental Setup

- ▶ Data: SWBD with 260h speech
- ▶ Acoustic models: Factorized TDNN with LFMMI objective (Kaldi)
- ▶ Neural LMs: 2-layer LSTM and a 6-layer Transformer (PyTorch); 2-layer LSTM (Kaldi RNNLM)

## Experimental Results

### Effect of Estimation Methods

Model	$\epsilon$	Average	Weighted Average	Semi-Viterbi
Transformer	0.5	10.7	10.7	<b>10.6</b>
	0.05	10.6	10.6	<b>10.5</b>

Table 1: WERs on Hub5’00 (full set) of SWBD from the non-iterative lattice rescoring strategy with three estimation methods.

- ▶ Similar observation is observed with the LSTM LM

### Analysis of Iterative Rescoring

Rescoring Method	Hub5’00	Swb	Callhm
Score replacement	10.8	6.8	14.6
Non-iterative ( $\epsilon = 0.5$ )	10.6	6.8	14.3
Score replacement + Non-iterative	<b>10.3</b>	<b>6.6</b>	<b>14.0</b>

Table 2: WERs from proposed lattice rescoring strategies with a Transformer LM.

- ▶ The better performance of non-iterative rescoring compared with score replacement alone indicates the value of lattice expansion

## Experimental Results (cont.)

### Comparison with $n$ -gram Expansion

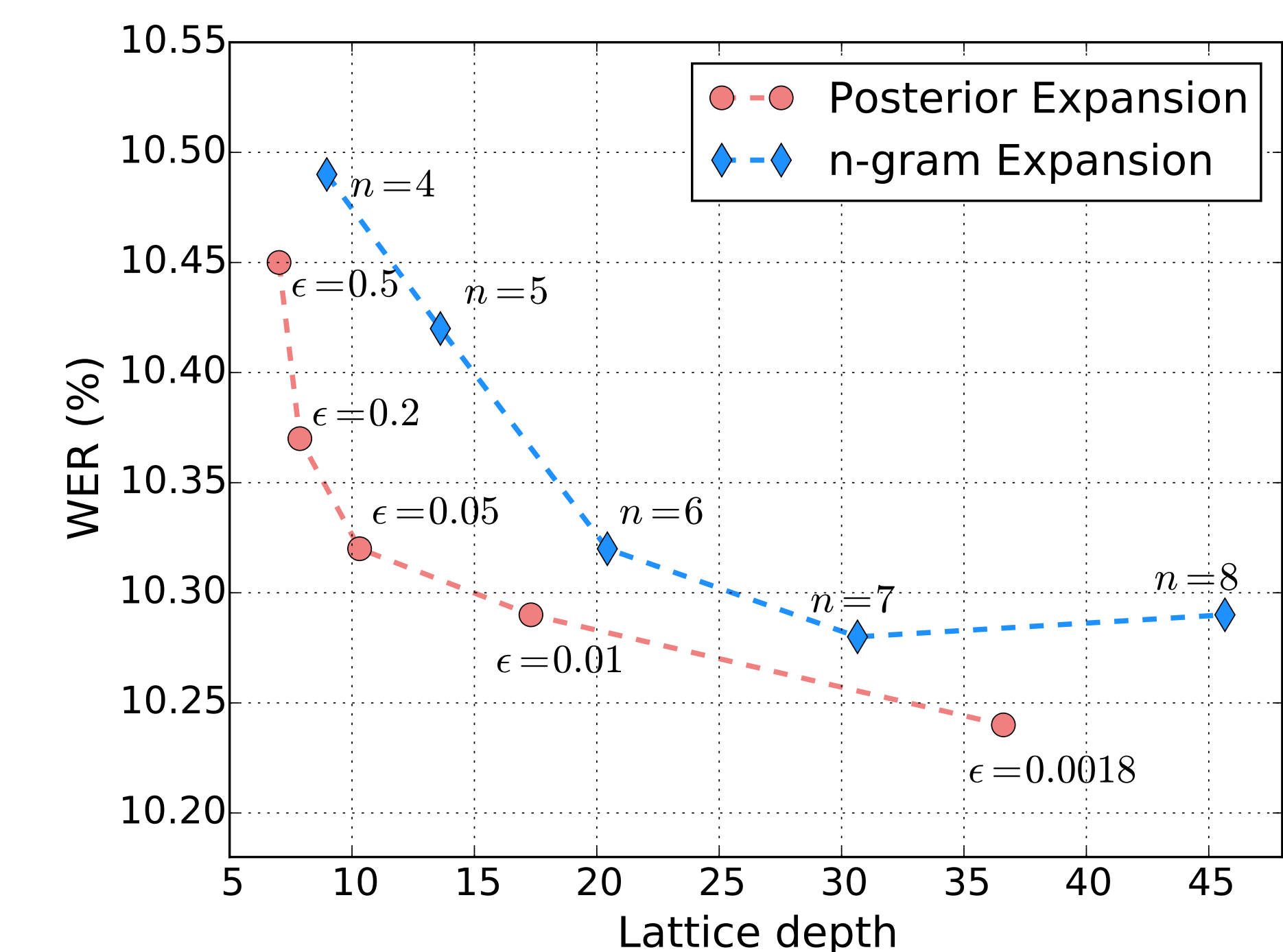


Figure 1: WERs and lattice depths for different  $\epsilon$  values and  $n$ -gram orders.

- ▶ Posterior-based expansion generates more compact lattices with better recognition accuracy than  $n$ -gram expansion

### Comparison with Pruned Lattice Rescoring

Method	WER			Lattice Depth
	Hub5’00	Swb	Callhm	
20-best	11.3	7.5	15.0	-
Pruned (4-gram approx.)	11.2	<b>7.3</b>	15.0	15.1
Non-iterative ( $\epsilon = 0.5$ )	<b>11.1</b>	7.4	<b>14.9</b>	<b>6.4</b>

Table 3: WERs and lattice depths from pruned lattice rescoring and the proposed non-iterative lattice rescoring.

- ▶ The non-iterative rescoring strategy obtains competitive performance and generates smaller lattices

### WERs on SWBD

Method	Hub5’00	Swb	Callhm
4-gram KN	12.8	8.6	17.0
20-best (LSTM)	10.9	7.1	14.6
20-best (Transformer)	10.8	7.2	14.4
Non-iterative ( $\epsilon = 0.005$ )	10.4	6.8	14.0
Iterative ( $\epsilon = 0.001$ )	<b>10.2</b>	<b>6.5</b>	<b>13.9</b>

Table 4: WERs from proposed lattice rescoring strategies with a Transformer LM.

## Conclusions

- ▶ Lattice-to-list conversion enables parallel LM evaluations within a lattice and fully takes advantage of the parallel computation cross words of Transformer LMs for speedup.
- ▶ Posterior-based lattice expansion outperforms  $n$ -gram expansion.
- ▶ The proposed rescoring strategy makes it easier and more flexible to perform lattice rescoring with PyTorch LMs in Kaldi `kaldi/egs/swbd/s5c/local/pytorchnn/`