



# Curriculum Optimization for Low-resource Speech Recognition

Anastasia Kuznetsova<sup>†</sup>, Anurag Kumar<sup>\*†</sup>, Jennifer Drexler Fox<sup>\*</sup>, Francis Tyers<sup>†</sup>

<sup>†</sup>Indiana University Bloomington <sup>\*</sup>Rev.com, USA

## Background

- ▶ Resource constrained ASR until now remains a challenging task. Our *curriculum learning* based approach mitigates the lack of labelled training data;
- ▶ A ranking function is applied to input data and then the data is split into  $K$  equal portions (tasks). The set of tasks is called *curriculum*.
- ▶ We use the ranking function as the prior on the input data, as well as the learner's progress to optimize the sequence of ASR inputs with the help of *multi-armed bandit* (MAB) algorithms.

## Multi-Armed Bandit

- ▶ MAB is the concept from reinforcement learning (RL) domain. Here we describe its main components;
- ▶ The *agent* takes actions  $a_k \in \mathcal{A}_K$  inside the environment; it selects the best action following the *policy*  $\pi$ ;
- ▶ Policy  $\pi$  is a value function over all possible actions  $\mathcal{A}_K$ ;
- ▶ The environment generates the *reward*  $r$  reflecting the optimality of the current action;
- ▶ The agent adjusts the policy based on the reward to improve its actions in the future.

## Complexity metrics

- ▶ We hypothesize that the signal-based features have more significance for ASR than textual features and score input audios with *compression ratio* (CR) metric, where  $Size_{before}$  is the size of the audio before the compression and  $Size_{after}$  after compression:

$$CR = 1 - \frac{Size_{before}}{Size_{after}} \quad (1)$$

- ▶ Audios containing less noise compress more and have higher CR, on the other hand noisy audios have lower CR;
- ▶ We compare CR to text-based *sentence length* (SL) and *sentence norm* (SN) derived from the sentence embeddings;

## References

Graves et al., "Automated curriculum learning for neural networks," in Proc. ICML'17.

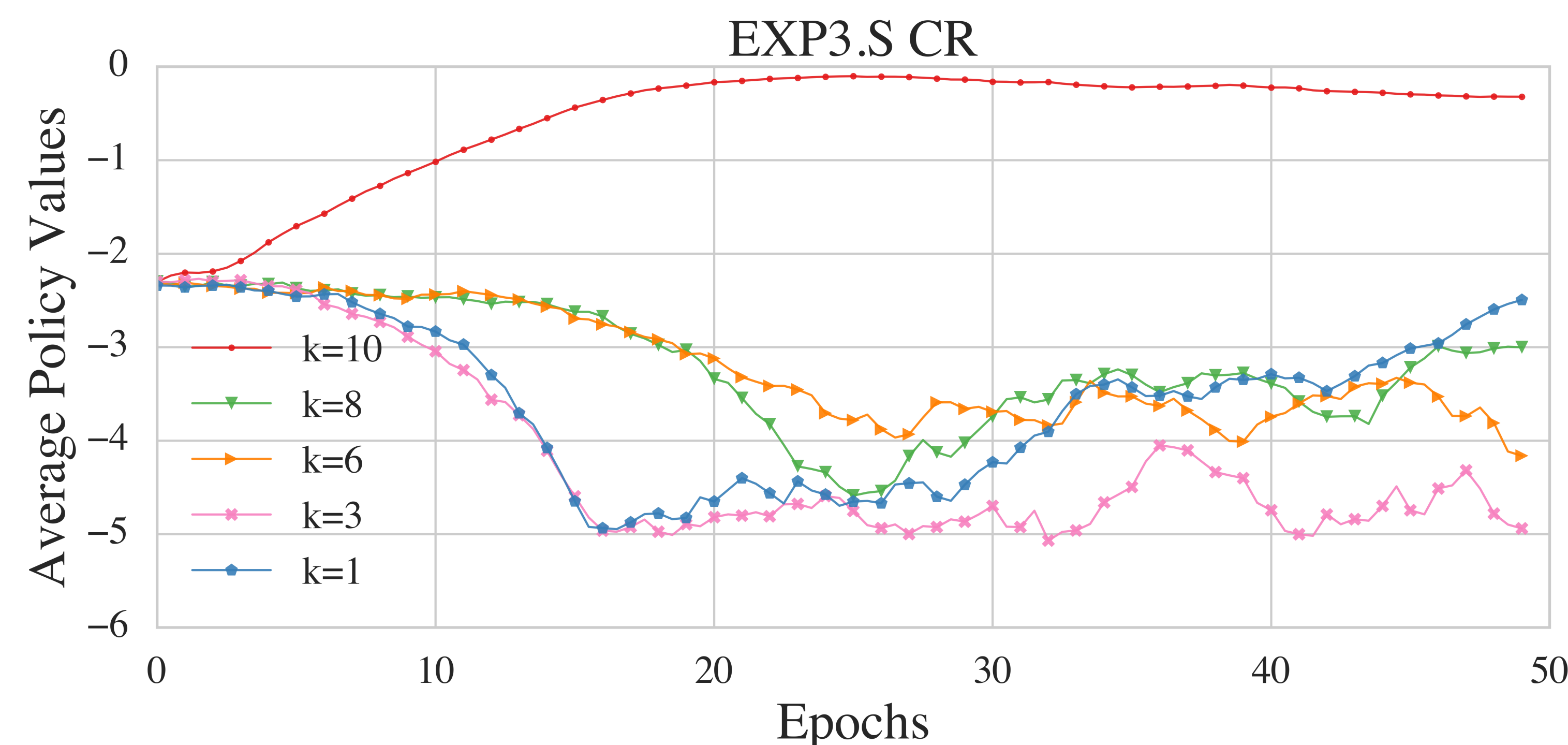


Figure: The figure shows log policy value change over time for selected  $k$ . This policy is generated by EXP3.S + CR on Kyrgyz (Ky) dataset and the policy is averaged per epoch.

Model	WER					CER				
	Cv	Fy	Tt	Ky	Eu	Cv	Fy	Tt	Ky	Eu
ESPNet + Trans	61.4	9.6	23.5	5.5	7.5	15.9	3.1	5.9	2.3	1.5
EXP3.S + CR	<b>41.8</b>	<b>7.8</b>	<b>22.3</b>	<b>4.1</b>	7.5	<b>9.8</b>	<b>2.6</b>	<b>5.6</b>	<b>1.9</b>	1.5
EXP3.S + SL	<b>41.1</b>	<b>9.3</b>	27.9	6.4	9.3	<b>9.7</b>	<b>3.3</b>	7.4	3.0	2.0
EXP3.S + SN	<b>42.7</b>	<b>7.8</b>	24.2	<b>4.8</b>	8.1	<b>10.2</b>	<b>2.5</b>	6.2	<b>2.2</b>	1.7
SW-UCB# + CR	<b>42.7</b>	<b>7.5</b>	<b>22.1</b>	<b>4.0</b>	8.0	<b>10.0</b>	<b>2.5</b>	<b>5.4</b>	<b>1.8</b>	1.6
SW-UCB# + SL	<b>42.4</b>	10.9	24.0	6.7	10.0	<b>9.9</b>	3.6	6.3	3.0	2.2
SW-UCB# + SN	<b>42.6</b>	<b>8.6</b>	24.8	<b>5.3</b>	8.6	<b>10.1</b>	<b>2.8</b>	6.4	<b>2.2</b>	1.8

Table: WER and CER for 5 selected languages. Baseline results are shown in the first row. Other rows show the combinations of the MAB algorithm and complexity metric, CR – Compression Ratio, SL – Sentence Length, SN – Sentence Norm for  $K = 10$ .

## Experimental setup

- ▶ The ASR model is trained on 5 Common Voice 7.0 languages using ESPNet recipe; 80-dimensional log-mel filterbanks are extracted from the audio, with a window length of 25ms and a stride of 10ms; The data is augmented with SpecAugment;
- ▶ The encoder has 9 Conformer blocks with 4 attention heads. The decoder has 6 transformer blocks. The network is trained with Adam optimizer with 25k steps warmup;
- ▶ The model is pretrained on Common Voice 7.0 English data set with the resulting WER of 15.2;
- ▶ The pretrained and fine-tuned on a target language models are used as a transfer learning baseline.

## Method

- ▶ The input audios are ranked and split into  $K$  tasks with equal number of mini-batches;
- ▶ At each iteration  $t$  the MAB selects the best action  $k$  and updates the policy based on the reward  $r_t$ ;
- ▶ The  $r_t$  is calculated from the loss-driven *self-prediction gain* (SPG)

$$\nu_{SPG} = L(\mathcal{B}', \theta) - L(\mathcal{B}', \theta') \quad \mathcal{B}' \sim D_k \quad (2)$$

where  $\mathcal{B}'$  is the batch sampled from task  $D_k$  and  $\theta$  are the parameters of the neural ASR model (Graves et al., 2017);

- ▶ We experiment with two MAB algorithms, probabilistic EXP3.S and deterministic SW-UCB# adapted for non-stationary problems.

## Algorithm 1: Curriculum Learning

**Initialize:**  $D = f(X)$ ,  $\pi \leftarrow 0$ ;

**begin**

**for**  $t \rightarrow T$  **do**

    Draw  $k$  based on current  $\pi$ ;

$\mathcal{B}_{t,k} \leftarrow \text{sample}(D_k)$ ;

    Train the model on  $\mathcal{B}_{t,k}$ ;

    Observe progress gain  $\nu_{SPG}$ ;

$r_t \leftarrow g(\nu_{SPG})$ ;

    Update  $\pi$  on  $r_t$ ;

**end**

**end**

## Results

- ▶ Our system achieved the highest relative decrease in WER of 33% relative for Chuvash (CV) and the minimum improvement of 5% for Tatar (Tt) language. For Fy, Tt and Ky, the best results are delivered with CR. For Cv the best WER was achieved by SL, however, the second best result is still achieved by CR.
- ▶ The policy over time (see Fig.) shows that harder  $k = 10$  is preferred earlier in the training and policy value for easy task  $k = 1$  increases towards the end of the training since the model gets more information from the harder examples in the beginning.