# Curriculum Optimization for Low Resource Speech Recognition

A. Kuznetsova[†], A. Kumar[★], J. Drexler Fox[★], F.Tyers[†]

[†]Indiana University Bloomington, [★]Rev.com, USA

INDIANA UNIVERSITY
**LUDDY SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING**

# Introduction

- Despite the recent achievements in speech recognition (ASR) domain, resource constrained ASR is still a challenge;
- Proposed approach is for labelled ASR and accounts for a lack of training data and diverse noise/quality conditions based on the idea of **curriculum learning**.

**Curriculum learning** is a machine learning strategy of imitating human study behavior;

- We define **curriculum** as a set of tasks organized in order of increasing complexity. A ranking function is used to determine the complexity of the input data;
- The proposed method aims to leverage the prior distribution of the training data as well as the learner's progress to optimize the sequence of ASR inputs by using **multi-armed bandit** (MAB) [2,3].

# Contributions

- We proposed a new complexity metric, *compression ratio,* able to rank audio signals in diverse noise/background conditions;

- Our model achieves the maximum of 33% and the minimum of 5% relative WER improvement;

- We show the analysis of the curriculum policy generated by MAB algorithms.

INDIANA UNIVERSITY
**LUDDY SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING**

# Curriculum Learning & Complexity Measures

In curriculum learning the model is fed input data scored by a ranking function in order of increasing complexity; we **hypothesize** that signal-based features have more significance for ASR than text-based features.
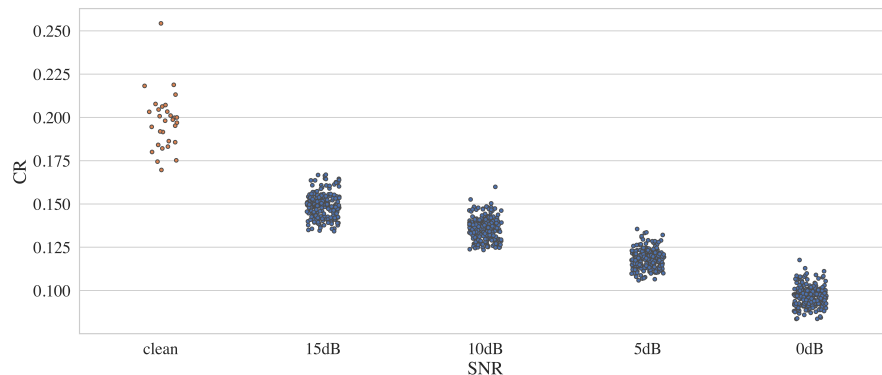
**Compression ratio (CR)**

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

**Sentence Norm (SN)**

$$SN = \left\| \frac{1}{N} \sum_{j=1}^{N} (y_j) \right\|_2$$

**Sentence length (SL)**



Compression ratio of clean audio compared to noisy mixtures at different SNRs (data from NOIZEUS) [4].

# Curriculum Definition

- Given a ranking function $f(\cdot)$ the inputs are scored and sorted in non-increasing order and split into a set of $K$ tasks with equal number of mini-batches:
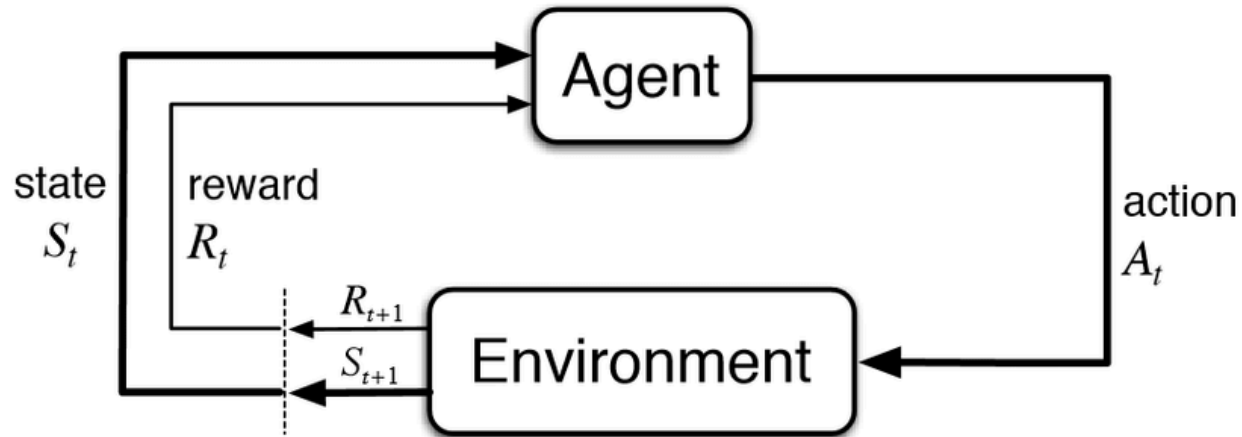
$$D = \{D_1, D_2, D_3, \ldots, D_K\}$$

- To use the prior distribution of the training data and the learner's progress simultaneously we apply multi-armed bandit algorithms.

# Multi-Armed Bandits (MAB)

Reinforcement Learning (RL) concepts: *agent, environment, policy $\pi$, reward $R$.*



Source: Sutton and Barto, Reinforcement Learning: Introduction, 1998.

# Multi-Armed Bandits (MAB) Cont.

Our **agent** is represented by the $k$-armed bandit which aims to collect the maximum expected reward $R$ over $T$ training iterations. At each iteration $t$ the agent selects the best action $k$ and updates the policy based on the reward

The reward based on the loss-driven *self-prediction gain* proposed in [1]:

$$\nu_{\mathrm{SPG}} = L(B', \theta) - L(B', \theta') \qquad B' \sim D_k$$

where $B'$ is the batch sampled from $D_k$;

$\theta$ denotes parameters of the model;

$L$ is the loss function.

**Probabilistic**: EXP3.S [2]

**Deterministic:** SWUCB# [3]

# Curriculum Learning Algorithm

**Algorithm 1:** Curriculum Learning

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

**begin**

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

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

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

        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**

# Experimental Setup

**ESPNet Common Voice Recipe [4]**

| Features | 80-dim log-mel filterbanks |
|---|---|
| **Augmentation** | SpecAugment |
| **Encoder** Conformer | 9 blocks 4 self-attention heads |
| **Decoder** Transformer | 6 blocks 4 self-attention heads |
| **Optimizer** | Adam , 25k steps warmup |
| **Pretraining** | CV 7.0 English, WER 15.2 |

**Common Voice 7.0 dataset**

|  | Eu | Fy | Ky | Tt | Cv |
|---|---|---|---|---|---|
| **Train** | 45:08 | 18:25 | 24:17 | 19:27 | 1:41 |
| **Dev** | 8:03 | 4:14 | 2:07 | 2:47 | 0:45 |
| **Test** | 8:34 | 4:25 | 2:12 | 5:07 | 1:08 |

# Results

| | WER | | | | | CER | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Model** | **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 | **41.8** | **7.8** | **22.3** | **4.1** | 7.5 | **9.8** | **2.6** | **5.6** | **1.9** | 1.5 |
| EXP3.S + SL | <u>41.1</u> | 9.3 | 27.9 | 6.4 | 9.3 | <u>9.7</u> | 3.3 | 7.4 | 3.0 | 2.0 |
| EXP3.S + SN | 42.7 | **7.8** | 24.2 | **4.8** | 8.1 | 10.2 | <u>2.5</u> | 6.2 | **2.2** | 1.7 |
| SW-UCB# + CR | 42.7 | <u>7.5</u> | <u>22.1</u> | **4.0** | 8.0 | **10.0** | <u>2.5</u> | <u>5.4</u> | <u>1.8</u> | 1.6 |
| SW-UCB# + SL | 42.4 | 10.9 | 24.0 | 6.7 | 10.0 | **9.9** | 3.6 | 6.3 | 3.0 | 2.2 |
| SW-UCB# + SN | **42.6** | **8.6** | 24.8 | **5.3** | 8.6 | **10.1** | 2.8 | 6.4 | **2.2** | 1.8 |

**Table 2**. The table shows WER and CER for five selected languages. Results in bold indicate the improvement over the baseline, underlined values indicate best result overall. Baseline results with bare transfer learning are shown in the first row. The results below show the combinations of the algorithm and complexity metric, $CR$ – Compression Ratio, $SL$ – Sentence Length, $SN$ – Sentence Norm for $K = 10$.

# Results Cont.



Policy values for all $k$ started with uniform distribution. The hardest task $k = 10$ initially gets the highest value; by the end of the training $k = 1$ increases the value.

# Conclusion

- The proposed method improves WER on 4 Common Voice languages: Cv, Fy, Tt, Ky;

- The combination of data derived prior and learner's progress curriculum yields better results;

- We confirmed out hypothesis that signal-based prior is more effective for ASR than text-based prior.

# References

1. *Graves, Alex, et al*. "Automated curriculum learning for neural networks." *international conference on machine learning*. PMLR, 2017.

2. *Auer, Peter, et al*. "The nonstochastic multiarmed bandit problem." *SIAM journal on computing* 32.1 (2002): 48-77.

3. *Wei, Lai, and Vaibhav Srivatsva*. "On abruptly-changing and slowly-varying multiarmed bandit problems." *2018 Annual American Control Conference (ACC)*. IEEE, 2018.

4. Watanabe, Shinji et al, "ESPnet: End-to-end speech processing toolkit," Proceedings of Interspeech, 2018, pp. 2207–2211.

INDIANA UNIVERSITY
**LUDDY SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING**

# Thank You!

Anastasia Kuznetsova
Department of Computer Science
Indiana University Bloomington
anakuzne@iu.edu