

---

# LeRaC: Learning Rate Curriculum

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Most curriculum learning methods require an approach to sort the data samples  
2 by difficulty, which is often cumbersome to perform. In this work, we propose a  
3 novel curriculum learning approach termed **Learning Rate Curriculum (LeRaC)**,  
4 which leverages the use of a different learning rate for each layer of a neural  
5 network to create a data-free curriculum during the initial training epochs. More  
6 specifically, LeRaC assigns higher learning rates to neural layers closer to the input,  
7 gradually decreasing the learning rates as the layers are placed farther away from  
8 the input. The learning rates increase at various paces during the first training  
9 iterations, until they all reach the same value. From this point on, the neural model  
10 is trained as usual. This creates a model-level curriculum learning strategy that  
11 does not require sorting the examples by difficulty and is compatible with any  
12 neural network, generating higher performance levels regardless of the architecture.  
13 We conduct comprehensive experiments on eight datasets from the computer vision  
14 (CIFAR-10, CIFAR-100, Tiny ImageNet), language (BoolQ, QNLI, RTE) and  
15 audio (ESC-50, CREMA-D) domains, considering various convolutional (ResNet-  
16 18, Wide-ResNet-50, DenseNet-121), recurrent (LSTM) and transformer (CvT,  
17 BERT, SepTr) architectures, comparing our approach with the conventional training  
18 regime. Moreover, we also compare with Curriculum by Smoothing (CBS), a state-  
19 of-the-art data-free curriculum learning approach. Unlike CBS, our performance  
20 improvements over the standard training regime are consistent across all datasets  
21 and models. Furthermore, we significantly surpass CBS in terms of training time  
22 (there is no additional cost over the standard training regime for LeRaC). Our code  
23 is freely available at: <http://github.com/link.hidden.for.review>.

## 24 1 Introduction

25 Machine learning researchers relentlessly strive to improve the performance of AI models. Much of  
26 this effort has been directed to the development of novel neural architectures [1–9], which have grown  
27 in size and complexity [1, 7, 10] to leverage the availability of increasingly larger datasets. However,  
28 we believe the dominant trend to develop deeper and deeper neural networks is not sustainable on  
29 the long term. To this end, we turn our attention to an alternative approach to increase performance  
30 of deep neural models without growing the size of the respective models. More specifically, we  
31 focus on curriculum learning, an approach initially proposed by Bengio et al. [11] to train better  
32 neural networks by mimicking how humans learn, from easy to hard. As originally introduced by  
33 Bengio et al. [11], curriculum learning is a training procedure that first organizes the examples in their  
34 increasing order of difficulty, then starts the training of the neural network on the easiest examples,  
35 gradually adding increasingly more difficult examples along the way, until all training examples  
36 are fed to the network. The success of the approach relies in avoiding to force the learning of very  
37 difficult examples right from the beginning, instead guiding the model on the right path through the  
38 imposed curriculum. This type of curriculum is later referred to as data-level curriculum learning  
39 [12]. Indeed, Soviany et al. [12] identified several types of curriculum learning approaches in the

40 literature, dividing them into four categories based on the components involved in the definition of  
41 machine learning given by Mitchell [13]. The four categories are: data-level curriculum (examples  
42 are presented from easy to hard), model-level curriculum (the modeling capacity of the network is  
43 gradually increased), task-level curriculum (the complexity of the learning task is increased during  
44 training), objective-level curriculum (the model optimizes towards an increasingly more complex  
45 objective). While data-level curriculum is the most natural and direct way to employ curriculum  
46 learning, its main disadvantage is that it requires a way to determine the difficulty of the data samples.  
47 The task of estimating the difficulty of the data samples has been addressed in different domain-  
48 specific ways, e.g. the length of text has been used in natural language processing [14, 15], while  
49 the number or size of objects were shown to work well in computer vision [16, 17]. Despite having  
50 many successful applications [12, 18], there is no universal way to determine the difficulty of the  
51 data samples, making the data-level curriculum less applicable to scenarios where the difficulty is  
52 hard to estimate, e.g. classification of radar signals. The task-level and objective-level curriculum  
53 learning strategies suffer from similar issues, e.g. it is hard to create a curriculum when the model has  
54 to learn an easy task (binary classification) or the objective function is already convex.

55 Considering the above observations, we recognize the potential of model-level curriculum learning  
56 strategies of being applicable across a wider range of domains and tasks. To date, there are only a few  
57 works [19–21] in the category of pure model-level curriculum learning methods. Furthermore, the  
58 existing methods have some drawbacks caused by their domain-dependent or architecture-specific  
59 design. For instance, Karras et al. [20] gradually increase the resolution of input images as new  
60 layers are appended to a generative network, but the notion of input resolution does not exist in other  
61 domains, e.g. text. Burduja et al. [19] blur the input images with Gaussian kernels, but this method is  
62 not applicable to an input format for which there is no convolution operation, e.g. tabular data. Sinha  
63 et al. [21] apply Gaussian kernel smoothing on convolutional activation maps, but this operation  
64 makes less sense for a feed-forward neural network formed only of dense layers.

65 To benefit from the full potential of the model-level curriculum learning category, we propose LeRaC  
66 (**L**earning **R**ate **C**urriculum), a novel and simple curriculum learning approach which leverages the  
67 use of a different learning rate for each layer of a neural network to create a data-free curriculum  
68 during the initial training epochs. More specifically, LeRaC assigns higher learning rates to neural  
69 layers closer to the input, gradually decreasing the learning rates as the layers are placed farther away  
70 from the input. This prevents the propagation of noise caused by the random initialization of the  
71 network’s weights. The learning rates increase at various paces during the first training iterations,  
72 until they all reach the same value. From this point on, the neural model is trained as usual. This  
73 creates a model-level curriculum learning strategy that is applicable to any domain and compatible  
74 with any neural network, generating higher performance levels regardless of the architecture, without  
75 adding any extra training time. To the best of our knowledge, we are the first to employ a different  
76 learning rate per layer to achieve the same effect as conventional (data-level) curriculum learning.

77 We conduct comprehensive experiments on eight datasets from the computer vision (CIFAR-10 [22],  
78 CIFAR-100 [22], Tiny ImageNet [23]), language (BoolQ [24], QNLI [25], RTE [25]) and audio (ESC-  
79 50 [26], CREMA-D [27]) domains, considering various convolutional (ResNet-18 [4], Wide-ResNet-  
80 50 [28], DenseNet-121 [29]), recurrent (LSTM [30]) and transformer (CvT [8], BERT [2], SepTr  
81 [31]) architectures, comparing our approach with the conventional training regime and Curriculum by  
82 Smoothing (CBS) [21], our closest competitor. Unlike CBS, our performance improvements over the  
83 standard training regime are consistent across all datasets and models. Furthermore, we significantly  
84 surpass CBS in terms of training time, since there is no additional cost over the conventional training  
85 regime for LeRaC, whereas CBS adds Gaussian kernel smoothing layers.

86 In summary, our contributions are twofold:

- 87 • We propose a novel and simple model-level curriculum learning strategy that creates a  
88 curriculum by updating the weights of each neural layer with a different learning rate,  
89 considering higher learning rates for the low-level feature layers and lower learning rates for  
90 the high-level feature layers.
- 91 • We empirically demonstrate the applicability to multiple domains (image, audio and text),  
92 the compatibility to several neural network architectures (convolutional neural networks,  
93 recurrent neural networks and transformers), and the time efficiency (no extra training time  
94 added) of LeRaC through a comprehensive set of experiments.

## 95 2 Related Work

96 Curriculum learning was initially introduced by Bengio et al. [11] as a training strategy that helps  
97 machine learning models to generalize better when the training examples are presented in the  
98 ascending order of their difficulty. Extensive surveys on curriculum learning methods, including  
99 the most recent advancements on the topic, were conducted by Soviany et al. [12] and Wang et  
100 al. [18]. In the former survey, Soviany et al. [12] emphasized that curriculum learning is not only  
101 applied at the data level, but also with respect to the other components involved in a machine  
102 learning approach, namely at the model level, the task level and the objective (performance measure)  
103 level. Regardless of the component on which curriculum learning is applied, the technique has  
104 demonstrated its effectiveness on a broad range of machine learning tasks, from computer vision  
105 [11, 16, 17, 21, 32–34] to natural language processing [11, 35–38] and audio processing [39, 40].

106 The main challenge for the methods that build the curriculum at the data level is measuring the  
107 difficulty of the data samples, which is required to order the samples from easy to hard. Most studies  
108 have addressed the problem with human input [41–43] or metrics based on domain-specific heuristics.  
109 For instance, the length of the sentence [36, 44] and the word frequency [11, 38] have been employed  
110 in natural language processing. In computer vision, the samples containing fewer and larger objects  
111 have been considered to be easier in some works [16, 17]. Other solutions employed difficulty  
112 estimators [45] or even the confidence level of the predictions made by the neural network [46, 47] to  
113 approximate the complexity of the data samples.

114 The solutions listed above have shown their utility in specific application domains. Nonetheless,  
115 measuring the difficulty remains problematic when implementing standard (data-level) curriculum  
116 learning strategies, at least in some application domains. Therefore, several alternatives have emerged  
117 over time, handling the drawback and improving the conventional curriculum learning approach. In  
118 [48], the authors introduced self-paced learning to evaluate the learning progress when selecting the  
119 easy samples. The method was successfully employed in multiple settings [48–54]. Furthermore,  
120 some studies combined self-paced learning with the traditional pre-computed difficulty metrics  
121 [53, 55]. An additional advancement related to self-paced learning is the approach called self-paced  
122 learning with diversity [56]. The authors demonstrated that enforcing a certain level of variety among  
123 the selected examples can improve the final performance. Another set of methods that bypass the  
124 need for predefined difficulty metrics is known as teacher-student curriculum learning [57, 58]. In  
125 this setting, a teacher network learns a curriculum to supervise a student neural network.

126 Closer to our work, a few methods [19–21] proposed to apply curriculum learning at the model level,  
127 by gradually increasing the learning capacity (complexity) of the neural architecture. Such curriculum  
128 learning strategies do not need to know the difficulty of the data samples, thus having a great potential  
129 to be useful in a broad range of tasks. For example, Karras et al. [20] proposed to gradually add  
130 layers to generative adversarial networks during training, while increasing the resolution of the input  
131 images at the same time. They are thus able to generate realistic high-resolution images. However,  
132 their approach is not applicable to every domain, since there is no notion of resolution for some  
133 input data types, e.g. text. Sinha et al. [21] presented a strategy that blurs the activation maps of  
134 the convolutional layers using Gaussian kernel layers, reducing the noisy information caused by the  
135 network initialization. The blur level is progressively reduced to zero by decreasing the standard  
136 deviation of the Gaussian kernels. With this mechanism, they obtain a training procedure that allows  
137 the neural network to see simple information at the start of the process and more intricate details  
138 towards the end. Curriculum by Smoothing (CBS) [21] was only shown to be useful for convolutional  
139 architectures applied in the image domain. Although we found that CBS is applicable to transformers  
140 by blurring the tokens, it is not necessarily applicable to any neural architecture, e.g. standard feed-  
141 forward neural networks. As an alternative to CBS, Burduja et al. [19] proposed to apply the same  
142 smoothing process on the input image instead of the activation maps. The method was applied with  
143 success in medical image alignment. However, this approach is not applicable to natural language  
144 input, as it is not clear how to apply the blurring operation on the input text.

145 Different from Burduja et al. [19] and Karras et al. [20], our approach is applicable to various  
146 domains, including but not limited to natural language processing, as demonstrated throughout our  
147 experiments. To the best of our knowledge, the only competing model-level curriculum method  
148 which is applicable to various domains is CBS [21]. Unlike CBS, LeRaC does not introduce new  
149 operations, such as smoothing with Gaussian kernels, during training. As such, our approach does  
150 not increase the training time with respect to the conventional training regime, as later shown in

151 the experiments. In summary, we consider that the simplicity of our approach comes with many  
 152 important advantages: applicability to any domain and task, compatibility with any neural network  
 153 architecture, time efficiency (adds no extra training time). We support all these claims through the  
 154 comprehensive experiments presented in Section 4.

### 155 3 Method

156 Deep neural networks are commonly trained on a set of labeled data samples denoted as:

$$S = \{(x_i, y_i) | x_i \in X, y_i \in Y, \forall i \in \{1, 2, \dots, m\}\}, \quad (1)$$

157 where  $m$  is the number of examples,  $x_i$  is a data sample and  $y_i$  is the associated label. The training  
 158 process of a neural network  $f$  with parameters  $\theta$  consists of minimizing some objective (loss) function  
 159  $\mathcal{L}$  that quantifies the differences between the ground-truth labels and the predictions of the model  $f$ :

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y_i, f(x_i, \theta)). \quad (2)$$

160 The optimization is generally performed by some variant of Stochastic Gradient Descent (SGD),  
 161 where the gradients are back-propagated from the neural layers closer to the output towards the neural  
 162 layers closer to input through the chain rule. Let  $f_1, f_2, \dots, f_n$  and  $\theta_1, \theta_2, \dots, \theta_n$  denote the neural  
 163 layers and the corresponding weights of the model  $f$ , such that the weights  $\theta_j$  belong to the layer  
 164  $f_j, \forall j \in \{1, 2, \dots, n\}$ . The output of the neural network for some training data sample  $x_i \in X$  is  
 165 formally computed as follows:

$$\hat{y}_i = f(x_i, \theta) = f_n(\dots f_2(f_1(x_i, \theta_1), \theta_2) \dots, \theta_n). \quad (3)$$

166 To optimize the model via SGD, the weights are updated as follows:

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta^{(t)} \cdot \frac{\partial \mathcal{L}}{\partial \theta_j^{(t)}}, \forall j \in \{1, 2, \dots, n\}, \quad (4)$$

167 where  $t$  is the index of the current training iteration,  $\eta^{(t)} > 0$  is the learning rate at iteration  $t$ , and the  
 168 gradient of  $\mathcal{L}$  with respect to  $\theta_j^{(t)}$  is computed via the chain rule. Before starting the training process,  
 169 the weights  $\theta_j^{(0)}$  are commonly initialized with random values.

170 Due to the random initialization of the weights, the information propagated through the neural model  
 171 during the early training iterations can contain a large amount of noise [21], which can negatively  
 172 impact the learning process. Due to the feed-forward processing, we conjecture that the noise level  
 173 tends to grow with each neural layer, from  $f_j$  to  $f_{j+1}$ . The same issue can occur if the weights are  
 174 pre-trained on a distinct task, where the misalignment of the weights with a new task is likely higher  
 175 for the high-level feature layers. To alleviate this problem, we propose to introduce a curriculum  
 176 learning strategy that assigns a different learning rate  $\eta_j$  to each layer  $f_j$ , as follows:

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \eta_j^{(t)} \cdot \frac{\partial \mathcal{L}}{\partial \theta_j^{(t)}}, \forall j \in \{1, 2, \dots, n\}, \quad (5)$$

177 such that:

$$\eta^{(0)} \geq \eta_1^{(0)} \geq \eta_2^{(0)} \geq \dots \geq \eta_n^{(0)}, \quad (6)$$

178

$$\eta^{(k)} = \eta_1^{(k)} = \eta_2^{(k)} = \dots = \eta_n^{(k)}, \quad (7)$$

179 where  $\eta_j^{(0)}$  are the initial learning rates and  $\eta_j^{(k)}$  are the updated learning rates at iteration  $k$ . The  
 180 condition formulated in Eq. (6) indicates that the initial learning rate  $\eta_j^{(0)}$  of a neural layer  $f_j$  gets  
 181 lower as the level of the respective neural layer becomes higher (farther away from the input). With  
 182 each training iteration  $t \leq k$ , the learning rates are gradually increased, until they become equal,  
 183 according to Eq. (7). Thus, our curriculum learning strategy is only applied during the early training  
 184 iterations, where the noise caused by the random weight initialization is most prevalent. Hence,  $k$  is a  
 185 hyperparameter of LeRaC that is usually adjusted such that  $k \ll T$ , where  $T$  is the total number of

186 training iterations. In practice, we obtain optimal results by running LeRaC up to any epoch between  
187 2 and 7.

188 We increase each learning rate  $\eta_j$  from iteration 0 to iteration  $k$  using an exponential scheduler that is  
189 based on the following rule:

$$\eta_j^{(l)} = \eta_j^{(0)} \cdot c^{\frac{l}{k} \cdot (\log_c \eta_j^{(k)} - \log_c \eta_j^{(0)})}, \forall l \in \{0, 1, \dots, k\}. \quad (8)$$

190 We set  $c = 10$  in Eq. (8) across all our experiments. In practice, we obtain optimal results by  
191 initializing the lowest learning rate  $\eta_n^{(0)}$  with a value that is around five or six orders of magnitude  
192 lower than  $\eta^{(0)}$ , while the highest learning rate  $\eta_1^{(0)}$  is usually equal to  $\eta^{(0)}$ . Apart from these general  
193 practical notes, the exact LeRaC configuration for each neural architecture is established by tuning  
194 the hyperparameters on the available validation sets.

## 195 4 Experiments

### 196 4.1 Datasets

197 In general, we adopt the official data splits for the eight benchmarks considered in our experiments.  
198 When a validation set is not available, we keep 10% of the training data for validation.

199 **CIFAR-10.** CIFAR-10 [22] is a popular dataset for object recognition in images. It consists of 60,000  
200 color images with a resolution of  $32 \times 32$  pixels. An images depicts one of 10 object classes, each  
201 class having 6,000 examples. We use the official data split with a training set of 50,000 images and a  
202 test set of 10,000 images.

203 **CIFAR-100.** The CIFAR-100 [22] dataset is similar to CIFAR-10, except that it has 100 classes with  
204 600 images per class. There are 50,000 training images and 10,000 test images.

205 **Tiny ImageNet.** Tiny ImageNet is a subset of ImageNet [23] which provides 100,000 training images,  
206 25,000 validation images and 25,000 test images representing objects from 200 different classes. The  
207 size of each image is  $64 \times 64$  pixels.

208 **BoolQ.** BoolQ [24] is a question answering dataset for yes/no questions containing 15,942 examples.  
209 The questions are naturally occurring, being generated in unprompted and unconstrained settings.  
210 Each example is a triplet of the form: {question, passage, answer}. We use the data split provided in  
211 the SuperGLUE benchmark [59], containing 9,427 examples for training, 3,270 for validation and  
212 3,245 for testing.

213 **QNLI.** The QNLI (Question-answering NLI) dataset [25] is a natural language inference benchmark  
214 automatically derived from SQuAD [60]. The dataset contains {question, sentence} pairs and the  
215 task is to determine whether the context sentence contains the answer to the question. The dataset  
216 is constructed on top of Wikipedia documents, each document being accompanied, on average, by  
217 4 questions. We consider the data split provided in the GLUE benchmark [25], which comprises  
218 104,743 examples for training, 5,463 for validation and 5,463 for testing.

219 **RTE.** Recognizing Textual Entailment (RTE) [25] is a natural language inference dataset containing  
220 pairs of sentences with the target label indicating if the meaning of one sentence can be inferred from  
221 the other. The training subset includes 2,490 samples, the validation set 277, and the test set 3,000  
222 examples.

223 **CREMA-D.** The CREMA-D multi-modal database [27] is formed of 7,442 videos of 91 actors (48  
224 male and 43 female) of different ethnic groups. The actors perform various emotions while uttering  
225 12 particular sentences that evoke one of the 6 emotion categories: anger, disgust, fear, happy, neutral,  
226 and sad. Following [54], we conduct experiments only on the audio modality, dividing the set of  
227 audio samples into 70% for training, 15% for validation and 15% for testing.

228 **ESC-50.** The ESC-50 [26] dataset is a collection of 2,000 samples of 5 seconds each, comprising 50  
229 classes of various common sound events. Samples are recorded at a 44.1 kHz sampling frequency,  
230 with a single channel. In our evaluation, we employ the 5-fold cross-validation procedure, as described  
231 in related works [26, 31].

Table 1: Optimal hyperparameter settings for the various neural architectures used in our experiments.

Architecture	Optimizer	Mini-batch	#Epochs	$\eta^{(0)}$	CBS			LeRaC		
					$\sigma$	$d$	$u$	$k$	$\eta_1^{(0)}$	$\eta_n^{(0)}$
ResNet-18	SGD	64	100-200	$10^{-1}$	1	0.9	2-5	5-7	$10^{-1}$	$10^{-8}$
Wide-ResNet-50	SGD	64	100-200	$10^{-1}$	1	0.9	2-5	5-7	$10^{-1}$	$10^{-8}$
CvT-13	Adamax	64-128	150-200	$2 \cdot 10^{-3}$	1	0.9	2-5	2-5	$2 \cdot 10^{-3}$	$2 \cdot 10^{-8}$
CvT-13 <sub>pre-trained</sub>	Adamax	64-128	25	$5 \cdot 10^{-4}$	1	0.9	2-5	3-6	$5 \cdot 10^{-4}$	$5 \cdot 10^{-10}$
BERT <sub>large-uncased</sub>	Adamax	10	7-25	$5 \cdot 10^{-5}$	1	0.9	1	3	$5 \cdot 10^{-5}$	$5 \cdot 10^{-8}$
LSTM	AdamW	256-512	25-70	$10^{-3}$	1	0.9	2	3-4	$10^{-3}$	$10^{-7}$
SepTR	Adam	2	50	$10^{-4}$	0.8	0.9	1-3	2-5	$10^{-4}$	$10^{-8}$
DenseNet-121	Adam	64	50	$10^{-4}$	0.8	0.9	1-3	2-5	$10^{-4}$	$5 \cdot 10^{-8}$

232 **4.2 Experimental Setup**

233 **Architectures.** To demonstrate the compatibility of LeRaC with multiple neural architectures, we  
 234 select several convolutional, recurrent and transformer models. As representative convolutional  
 235 neural networks (CNNs), we opt for ResNet-18 [4], Wide-ResNet-50 [28] and DenseNet-121 [29].  
 236 As representative transformers, we consider CvT-13 [8], BERT<sub>uncased-large</sub> [2] and SepTr [31]. For  
 237 CvT, we consider both pre-trained and randomly initialized versions. We use an uncased large pre-  
 238 trained version of BERT. As Ristea et al. [31], we train SepTr from scratch. In addition, we employ  
 239 a long short-term memory (LSTM) network [30] to represent recurrent neural networks (RNNs).  
 240 The recurrent neural network contains two LSTM layers, each having a hidden dimension of 256  
 241 components. These layers are preceded by one embedding layer with the embedding size set to 128  
 242 elements. The output of the recurrent layers is passed to a classifier comprised of two fully connected  
 243 layers. The LSTM is activated by rectified linear units (ReLU). We apply the aforementioned models  
 244 on distinct input data types, considering the intended application domain of each model<sup>1</sup>. Hence,  
 245 ResNet-18, Wide-ResNet-50 and CvT are applied on images, BERT and LSTM are applied on text,  
 246 and SepTr and DenseNet-121 are applied on audio.

247 **Baselines.** We compare LeRaC with two baselines: the conventional training regime (which uses  
 248 early stopping and reduces the learning rate on plateau) and the state-of-the-art Curriculum by  
 249 Smoothing [21]. For CBS, we use the official code released by Sinha et al. [21] at <https://github.com/pairlab/CBS>, to ensure the replicability of their method in our experimental settings, which  
 250 include a more diverse selection of input data types and neural architectures.  
 251

252 **Hyperparameter tuning.** We tune all hyperparameters on the validation set of each benchmark.  
 253 In Table 1, we present the optimal hyperparameters chosen for each architecture. In addition to the  
 254 standard parameters of the training process, we report the parameters that are specific for the CBS  
 255 and LeRaC strategies. In the case of CBS,  $\sigma$  denotes the standard deviation of the Gaussian kernel,  $d$   
 256 is the decay rate for  $\sigma$ , and  $u$  is the decay step. Regarding the parameters of LeRaC,  $k$  represents  
 257 the number of iterations used in Eq. (8), and  $\eta_1^{(0)}$  and  $\eta_n^{(0)}$  are the initial learning rates for the first  
 258 and last layers of the architecture, respectively. We underline that  $\eta_1^{(0)} = \eta^{(0)}$  and  $c = 10$ , in all  
 259 experiments. Moreover,  $\eta_j^{(k)} = \eta^{(0)}$ , i.e. the initial learning rates of LeRaC converge to the original  
 260 learning rate set for the conventional training regime. All models are trained with early stopping and  
 261 the learning rate is reduced by a factor of 10 when the loss reaches a plateau.

262 **Evaluation.** We evaluate all models in terms of the classification accuracy. We repeat the training  
 263 process of each model for 5 times and report the average accuracy and the standard deviation.

264 **Image preprocessing.** For the image classification experiments, we apply the same data preprocessing  
 265 approach as Sinha et al. [21]. Hence, we normalize the images and maintain their original resolution,  
 266  $32 \times 32$  pixels for CIFAR-10 and CIFAR-100, and  $64 \times 64$  pixels for Tiny ImageNet. Similar to  
 267 Sinha et al. [21], we do not employ data augmentation.

<sup>1</sup>The only exception is DenseNet-121, which is applied on audio instead of image data.

Table 2: Average accuracy rates (in %) over 5 runs on CIFAR-10, CIFAR-100 and Tiny ImageNet for various neural models based on different training regimes: conventional, CBS [21] and LeRaC. The accuracy of the best training regime in each experiment is highlighted in bold.

Architecture	Training Regime	CIFAR-10	CIFAR-100	Tiny ImageNet
ResNet-18	conventional	89.20 $\pm$ 0.43	65.28 $\pm$ 0.16	57.41 $\pm$ 0.05
ResNet-18	CBS [21]	89.53 $\pm$ 0.22	<b>66.41</b> $\pm$ 0.21	55.49 $\pm$ 0.20
ResNet-18	LeRaC (ours)	<b>89.56</b> $\pm$ 0.16	66.02 $\pm$ 0.17	<b>57.86</b> $\pm$ 0.20
Wide-ResNet-50	conventional	91.22 $\pm$ 0.24	68.14 $\pm$ 0.16	55.97 $\pm$ 0.30
Wide-ResNet-50	CBS [21]	89.05 $\pm$ 1.00	65.73 $\pm$ 0.36	48.30 $\pm$ 1.53
Wide-ResNet-50	LeRaC (ours)	<b>91.58</b> $\pm$ 0.16	<b>69.38</b> $\pm$ 0.26	<b>56.48</b> $\pm$ 0.60
CvT-13	conventional	71.84 $\pm$ 0.37	41.87 $\pm$ 0.16	33.38 $\pm$ 0.27
CvT-13	CBS [21]	72.64 $\pm$ 0.29	<b>44.48</b> $\pm$ 0.40	33.56 $\pm$ 0.36
CvT-13	LeRaC (ours)	<b>72.90</b> $\pm$ 0.28	43.46 $\pm$ 0.18	<b>33.95</b> $\pm$ 0.28
CvT-13 <sub>pre-trained</sub>	conventional	93.56 $\pm$ 0.05	77.80 $\pm$ 0.16	70.71 $\pm$ 0.35
CvT-13 <sub>pre-trained</sub>	CBS [21]	85.85 $\pm$ 0.15	62.35 $\pm$ 0.48	68.41 $\pm$ 0.13
CvT-13 <sub>pre-trained</sub>	LeRaC (ours)	<b>94.15</b> $\pm$ 0.03	<b>78.93</b> $\pm$ 0.05	<b>71.34</b> $\pm$ 0.08

Table 3: Average accuracy rates (in %) over 5 runs on BoolQ, RTE and QNLI for BERT and LSTM based on different training regimes: conventional, CBS [21] and LeRaC. The accuracy of the best training regime in each experiment is highlighted in bold.

Architecture	Training Regime	BoolQ	RTE	QNLI
BERT <sub>large-uncased</sub>	conventional	74.12 $\pm$ 0.32	74.48 $\pm$ 1.36	92.13 $\pm$ 0.08
BERT <sub>large-uncased</sub>	CBS [21]	74.37 $\pm$ 1.11	74.97 $\pm$ 1.96	91.47 $\pm$ 0.22
BERT <sub>large-uncased</sub>	LeRaC (ours)	<b>75.55</b> $\pm$ 0.66	<b>75.81</b> $\pm$ 0.29	<b>92.45</b> $\pm$ 0.13
LSTM	conventional	64.40 $\pm$ 1.37	54.12 $\pm$ 1.60	59.42 $\pm$ 0.36
LSTM	CBS [21]	64.75 $\pm$ 1.54	54.03 $\pm$ 0.45	59.89 $\pm$ 0.38
LSTM	LeRaC (ours)	<b>65.80</b> $\pm$ 0.33	<b>55.71</b> $\pm$ 1.04	<b>59.98</b> $\pm$ 0.34

268 **Text preprocessing.** For the text classification experiments with BERT, we lowercase all words and  
 269 add the classification token ([CLS]) at the start of the input sequence. We add the separator token  
 270 ([SEP]) to delimit sentences. For the LSTM network, we lowercase all words and replace them with  
 271 indexes from vocabularies constructed from the training set. The input sequence length is limited to  
 272 512 tokens for BERT and 200 tokens for LSTM.

273 **Speech preprocessing.** We transform each audio sample into a time-frequency matrix by computing  
 274 the discrete Short Time Fourier Transform (STFT) with  $N_x$  FFT points, using a Hamming window of  
 275 length  $L$  and a hop size  $R$ . For CREMA-D, we first standardize all audio clips to a fixed dimension  
 276 of 4 seconds by padding or clipping the samples. Then, we apply the STFT with  $N_x = 1024$ ,  
 277  $R = 64$  and a window size of  $L = 512$ . For ESC-50, we keep the same values for  $N_x$  and  $L$ , but  
 278 we increase the hop size to  $R = 128$ . Next, for each STFT, we compute the square root of the  
 279 magnitude and map the values to 128 Mel bins. The result is converted to a logarithmic scale and  
 280 normalized to the interval  $[0, 1]$ . Furthermore, in all our speech classification experiments, we use the  
 281 following data augmentation methods: noise perturbation, time shifting, speed perturbation, mix-up  
 282 and SpecAugment [61]. The speech preprocessing steps are carried out following Ristea et al. [31].

### 283 4.3 Results

284 **Image classification.** In Table 2, we present the image classification results on CIFAR-10, CIFAR-  
 285 100 and Tiny ImageNet. On the one hand, there are two scenarios (ResNet-18 on CIFAR-100 and  
 286 CvT-13 on CIFAR-100) in which CBS provides the largest improvements over the conventional  
 287 regime, surpassing LeRaC in the respective cases. On the other hand, there are seven scenarios  
 288 where CBS degrades the accuracy with respect to the standard training regime. This shows that the  
 289 improvements attained by CBS are inconsistent across models and datasets. Unlike CBS, our strategy

Table 4: Average accuracy rates (in %) over 5 runs on CREMA-D and ESC-50 for SepTr and DenseNet-121 based on different training regimes: conventional, CBS [21] and LeRaC. The accuracy of the best training regime in each experiment is highlighted in bold.

Architecture	Training Regime	CREMA-D	ESC-50
SepTr	conventional	70.47 ± 0.67	91.13 ± 0.33
SepTr	CBS [21]	69.98 ± 0.71	91.15 ± 0.41
SepTr	LeRaC (ours)	<b>70.95 ± 0.56</b>	<b>91.58 ± 0.28</b>
DenseNet-121	conventional	67.21 ± 0.12	88.91 ± 0.11
DenseNet-121	CBS [21]	68.16 ± 0.19	88.76 ± 0.17
DenseNet-121	LeRaC (ours)	<b>68.99 ± 0.08</b>	<b>90.02 ± 0.10</b>

Table 5: Average accuracy rates (in %) over 5 runs on CIFAR-10, CIFAR-100 and Tiny ImageNet for CvT-13 based on different training regimes: conventional, CBS [21], LeRaC with linear update, LeRaC with exponential update (proposed), and a combination of CBS and LeRaC.

Architecture	Training Regime	CIFAR-10	CIFAR-100	Tiny ImageNet
CvT-13	conventional	71.84 ± 0.37	41.87 ± 0.16	33.38 ± 0.27
CvT-13	CBS [21]	72.64 ± 0.29	44.48 ± 0.40	33.56 ± 0.36
CvT-13	LeRaC (linear update)	72.49 ± 0.27	43.39 ± 0.14	33.86 ± 0.07
CvT-13	LeRaC (exponential update)	72.90 ± 0.28	43.46 ± 0.18	33.95 ± 0.28
CvT-13	CBS [21] + LeRaC	73.25 ± 0.19	44.90 ± 0.41	34.20 ± 0.61

290 surpasses the baseline regime in all twelve cases, thus being more consistent. In four of these cases,  
 291 the accuracy gains of LeRaC are higher than 1%. Moreover, LeRaC outperforms CBS in ten out of  
 292 twelve cases. We thus consider that LeRaC can be regarded as a better choice than CBS, bringing  
 293 consistent performance gains.

294 **Text classification.** In Table 3, we report the text classification results on BoolQ, RTE and QNLI.  
 295 Here, there are only two cases (BERT on QNLI and LSTM on RTE) where CBS leads to performance  
 296 drops compared to the conventional training regime. In all other cases, the improvements of CBS are  
 297 below 0.6%. Just as in the image classification experiments, LeRaC brings accuracy gains for each  
 298 and every model and dataset. In four out of six scenarios, the accuracy gains yielded by LeRaC are  
 299 higher than 1.3%. Once again, LeRaC proves to be the best and most consistent regime, generally  
 300 outperforming CBS by significant margins.

301 **Speech classification.** In Table 4, we present the results obtained on the audio data sets, namely  
 302 CREMA-D and ESC-50. We observe that the CBS strategy obtains lower results compared with  
 303 the baseline in two cases (SepTr on CREMA-D and DenseNet-121 on ESC-50), while our method  
 304 provides superior results for each and every case. By applying LeRaC on SepTr, we set a new  
 305 state-of-the-art accuracy level (70.95%) on the CREMA-D audio modality, surpassing the previous  
 306 state-of-the-art value attained by Ristea et al. [31] with SepTr alone. When applied on DenseNet-121,  
 307 LeRaC brings performance improvements higher than 1%, the highest improvement (1.78%) over  
 308 the baseline being attained on CREMA-D.

309 **Additional results.** An interesting aspect worth studying is to determine if putting the CBS and  
 310 LeRaC regimes together could bring further performance gains. Across all our experiments, we  
 311 identified a single model (CvT-13) for which both CBS and LeRaC bring accuracy gains on all  
 312 datasets (see Table 2). We thus consider this model to try out the combination of CBS and LeRaC.  
 313 The corresponding results are shown in Table 5. The reported results show that the combination brings  
 314 accuracy gains across all three datasets (CIFAR-10, CIFAR-100, Tiny ImageNet). We thus conclude  
 315 that the combination of curriculum learning regimes is worth a try, whenever the two independent  
 316 regimes boost performance.

317 Another important aspect is to establish if the exponential learning rate update proposed in Eq. (8) is  
 318 a good choice. To test this out, we keep the CvT-13 model and change the LeRaC regime to use a  
 319 linear update of the learning rate. We observe performance gains with both types of update rules,



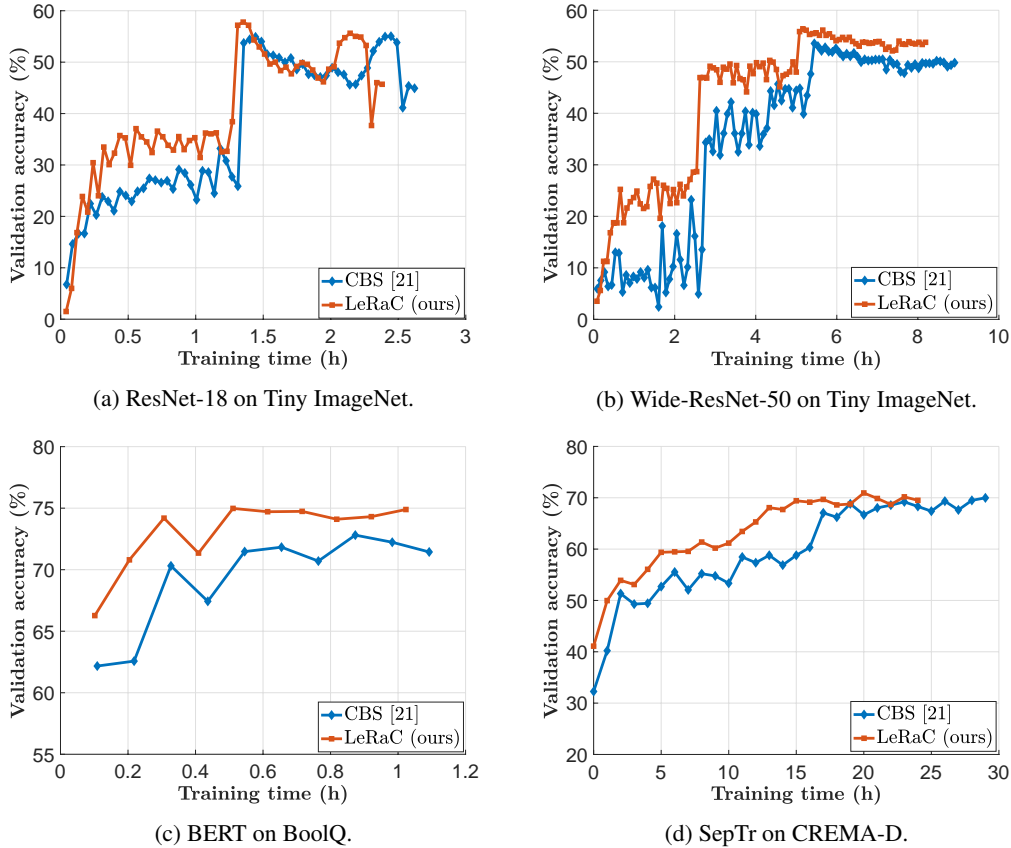


Figure 1: Validation accuracy (on the y-axis) versus training time (on the x-axis) for four distinct architectures. The number of training epochs is the same for both LeRaC and CBS, the observable time difference being caused by the overhead of CBS due to the Gaussian kernel layers.

320 but our exponential learning rate update seems to bring higher gains on all three datasets. We thus  
 321 conclude that the update rule defined in Eq. (8) is a sound option.

322 **Training time comparison.** For a particular model and dataset, all training regimes are executed for  
 323 the same number of epochs, for a fair comparison. However, the CBS strategy adds the smoothing  
 324 operation at multiple levels inside the architecture, which increases the training time. To this end,  
 325 we compare the training time (in hours) versus the validation error of CBS and LeRaC. For this  
 326 experiment, we selected four neural models and illustrate the evolution of the validation accuracy  
 327 over time in Figure 1. We underline that LeRaC introduces faster convergence times, being around  
 328 7-12% faster than CBS. It is trivial to note that LeRaC requires the same time as the conventional  
 329 regime.

## 330 5 Conclusion

331 In this paper, we introduced a novel model-level curriculum learning approach that is based on  
 332 starting the training process with increasingly lower learning rates per layer, as the layers get closer  
 333 to the output. We conducted comprehensive experiments on eight datasets from three domains  
 334 (image, text and audio), considering multiple neural architectures (CNNs, RNNs and transformers),  
 335 to compare our novel training regime (LeRaC) with a state-of-the-art regime (CBS [21]) as well as  
 336 the conventional training regime (based on early stopping and reduce on plateau). The empirical  
 337 results demonstrate that LeRaC is significantly more consistent than CBS, perhaps being the most  
 338 versatile curriculum learning strategy to date, due to its compatibility with multiple neural models  
 339 and its usefulness across different domains. Remarkably, all these benefits come for free, i.e. LeRaC  
 340 does not add any extra time over the conventional approach.

341 **References**

- 342 [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
343 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
344 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,  
345 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott  
346 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
347 Sutskever, and Dario Amodei, “Language Models are Few-Shot Learners,” in *Proceedings of*  
348 *NeurIPS*, 2020, vol. 33, pp. 1877–1901.
- 349 [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “BERT: Pre-training  
350 of Deep Bidirectional Transformers for Language Understanding,” in *Proceedings of NAACL*,  
351 2019, pp. 4171–4186.
- 352 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
353 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
354 Jakob Uszkoreit, and Neil Houlsby, “An Image is Worth 16x16 Words: Transformers for Image  
355 Recognition at Scale,” in *Proceedings of ICLR*, 2021.
- 356 [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep Residual Learning for Image  
357 Recognition,” in *Proceedings of CVPR*, 2016, pp. 770–778.
- 358 [5] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan,  
359 and Mubarak Shah, “Transformers in Vision: A Survey,” *ACM Computing Surveys*, 2021.
- 360 [6] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining  
361 Xie, “A ConvNet for the 2020s,” in *Proceedings of CVPR*, 2022.
- 362 [7] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
363 Yanqi Zhou, Wei Li, and Peter J. Liu, “Exploring the Limits of Transfer Learning with a Unified  
364 Text-to-Text Transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67,  
365 2020.
- 366 [8] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang,  
367 “CvT: Introducing Convolutions to Vision Transformers,” in *Proceedings of ICCV*, 2021, pp.  
368 22–31.
- 369 [9] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai, “Deformable DETR:  
370 Deformable Transformers for End-to-End Object Detection,” in *Proceedings of ICLR*, 2020.
- 371 [10] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov,  
372 Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom  
373 Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell,  
374 Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas, “A Generalist Agent,” *arXiv preprint*  
375 *arXiv:2205.06175*, 2022.
- 376 [11] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston, “Curriculum Learning,”  
377 in *Proceedings of ICML*, 2009, pp. 41–48.
- 378 [12] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe, “Curriculum learning: A  
379 survey,” *International Journal of Computer Vision*, 2022.
- 380 [13] Tom M. Mitchell, *Machine Learning*, McGraw-Hill, New York, 1997.
- 381 [14] Yi Tay, Shuohang Wang, Anh Tuan Luu, Jie Fu, Minh C. Phan, Xingdi Yuan, Jinfeng Rao,  
382 Siu Cheung Hui, and Aston Zhang, “Simple and Effective Curriculum Pointer-Generator  
383 Networks for Reading Comprehension over Long Narratives,” in *Proceedings of ACL*, 2019, pp.  
384 4922–4931.
- 385 [15] Wei Zhang, Wei Wei, Wen Wang, Lingling Jin, and Zheng Cao, “Reducing BERT Computation  
386 by Padding Removal and Curriculum Learning,” in *Proceedings of ISPASS*, 2021, pp. 90–92.
- 387 [16] Miaojing Shi and Vittorio Ferrari, “Weakly Supervised Object Localization Using Size Esti-  
388 mates,” in *Proceedings of ECCV*, 2016, pp. 105–121.

- 389 [17] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe, “Curriculum self-paced  
390 learning for cross-domain object detection,” *Computer Vision and Image Understanding*, vol.  
391 204, pp. 103–166, 2021.
- 392 [18] Xin Wang, Yudong Chen, and Wenwu Zhu, “A Survey on Curriculum Learning,” *IEEE*  
393 *Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- 394 [19] Mihail Burduja and Radu Tudor Ionescu, “Unsupervised Medical Image Alignment with  
395 Curriculum Learning,” in *Proceedings of ICIP*, 2021, pp. 3787–3791.
- 396 [20] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen, “Progressive Growing of GANs  
397 for Improved Quality, Stability, and Variation,” in *Proceedings of ICLR*, 2018.
- 398 [21] Samarth Sinha, Animesh Garg, and Hugo Larochelle, “Curriculum by smoothing,” in *Proceed-*  
399 *ings of NeurIPS*, 2020, pp. 21653–21664.
- 400 [22] Alex Krizhevsky, “Learning multiple layers of features from tiny images,” Tech. Rep., University  
401 of Toronto, 2009.
- 402 [23] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
403 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al., “ImageNet Large Scale  
404 Visual Recognition Challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp.  
405 211–252, 2015.
- 406 [24] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and  
407 Kristina Toutanova, “BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions,”  
408 in *Proceedings of NAACL*, 2019, pp. 2924–2936.
- 409 [25] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman,  
410 “GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding,”  
411 in *Proceedings of ICLR*, 2019.
- 412 [26] Karol J. Piczak, “ESC: Dataset for Environmental Sound Classification,” in *Proceedings of*  
413 *ACMMM*, 2015, pp. 1015–1018.
- 414 [27] Houwei Cao, David G. Cooper, Michael K. Keutmann, Ruben C. Gur, Ani Nenkova, and Ragini  
415 Verma, “CREMA-D: Crowd-sourced emotional multimodal actors dataset,” *IEEE Transactions*  
416 *on Affective Computing*, vol. 5, no. 4, pp. 377–390, 2014.
- 417 [28] Sergey Zagoruyko and Nikos Komodakis, “Wide Residual Networks,” *arXiv preprint*  
418 *arXiv:1605.07146*, 2016.
- 419 [29] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger, “Densely  
420 Connected Convolutional Networks,” in *Proceedings of CVPR*, 2017, pp. 2261–2269.
- 421 [30] Sepp Hochreiter and Jürgen Schmidhuber, “Long Short-Term Memory,” *Neural Computing*,  
422 vol. 9, no. 8, pp. 1735–1780, 1997.
- 423 [31] Nicolae-Catalin Ristea, Radu Tudor Ionescu, and Fahad Shahbaz Khan, “SepTr: Separable  
424 Transformer for Audio Spectrogram Processing,” *arXiv preprint arXiv:2203.09581*, 2022.
- 425 [32] Liangke Gui, Tadas Baltrušaitis, and Louis-Philippe Morency, “Curriculum Learning for Facial  
426 Expression Recognition,” in *Proceedings of FG*, 2017, pp. 505–511.
- 427 [33] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei, “MentorNet: Learning  
428 Data-Driven Curriculum for Very Deep Neural Networks on Corrupted Labels,” in *Proceedings*  
429 *of ICML*, 2018, pp. 2304–2313.
- 430 [34] Xinlei Chen and Abhinav Gupta, “Webly Supervised Learning of Convolutional Networks,” in  
431 *Proceedings of ICCV*, 2015, pp. 1431–1439.
- 432 [35] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabas Poczos, and Tom  
433 Mitchell, “Competence-based curriculum learning for neural machine translation,” in *Proceed-*  
434 *ings of NAACL*, 2019, pp. 1162–1172.

- 435 [36] Tom Kocmi and Ondřej Bojar, “Curriculum Learning and Minibatch Bucketing in Neural  
436 Machine Translation,” in *Proceedings of RANLP*, 2017, pp. 379–386.
- 437 [37] Valentin I. Spitzkovsky, Hiyan Alshawi, and Daniel Jurafsky, “Baby steps: How “less is more”  
438 in unsupervised dependency parsing,” in *Proceedings of NIPS*, 2009.
- 439 [38] Cao Liu, Shizhu He, Kang Liu, and Jun Zhao, “Curriculum Learning for Natural Answer  
440 Generation,” in *Proceedings of IJCAI*, 2018, pp. 4223–4229.
- 441 [39] Shivesh Ranjan and John H. L. Hansen, “Curriculum Learning Based Approaches for Noise  
442 Robust Speaker Recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language  
443 Processing*, vol. 26, pp. 197–210, 2018.
- 444 [40] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg,  
445 Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, Jie Chen, Jingdong  
446 Chen, Zhijie Chen, Mike Chrzanowski, Adam Coates, Greg Diamos, Ke Ding, Niandong Du,  
447 Erich Elsen, Jesse Engel, Weiwei Fang, Linxi Fan, Christopher Fougner, Liang Gao, Caixia  
448 Gong, Awni Hannun, Tony Han, Lappi Vaino Johannes, Bing Jiang, Cai Ju, Billy Jun, Patrick  
449 LeGresley, Libby Lin, Junjie Liu, Yang Liu, Weigao Li, Xiangang Li, Dongpeng Ma, Sharan  
450 Narang, Andrew Ng, Sherjil Ozair, Yiping Peng, Ryan Prenger, Sheng Qian, Zongfeng Quan,  
451 Jonathan Raiman, Vinay Rao, Sanjeev Satheesh, David Seetapun, Shubho Sengupta, Kavya  
452 Srinet, Anuroop Sriram, Haiyuan Tang, Liliang Tang, Chong Wang, Jidong Wang, Kaifu Wang,  
453 Yi Wang, Zhijian Wang, Zhiqian Wang, Shuang Wu, Likai Wei, Bo Xiao, Wen Xie, Yan Xie,  
454 Dani Yogatama, Bin Yuan, Jun Zhan, and Zhenyao Zhu, “Deep Speech 2: End-to-End Speech  
455 Recognition in English and Mandarin,” in *Proceedings of ICML*, 2016, pp. 173–182.
- 456 [41] Anastasia Pentina, Viktoriia Sharmanska, and Christoph H. Lampert, “Curriculum learning of  
457 multiple tasks,” in *Proceedings of CVPR*, June 2015, pp. 5492–5500.
- 458 [42] Amelia Jiménez-Sánchez, Diana Mateus, Sonja Kirchhoff, Chlodwig Kirchhoff, Peter Bib-  
459 erthaler, Nassir Navab, Miguel A. González Ballester, and Gemma Piella, “Medical-based Deep  
460 Curriculum Learning for Improved Fracture Classification,” in *Proceedings of MICCAI*, 2019,  
461 pp. 694–702.
- 462 [43] Jerry Wei, Arief Suriawinata, Bing Ren, Xiaoying Liu, Mikhail Lisovsky, Louis Vaickus, Charles  
463 Brown, Michael Baker, Mustafa Nasir-Moin, Naofumi Tomita, Lorenzo Torresani, Jason Wei,  
464 and Saeed Hassanpour, “Learn like a Pathologist: Curriculum Learning by Annotator Agreement  
465 for Histopathology Image Classification,” in *Proceedings of WACV*, 2021, pp. 2472–2482.
- 466 [44] Volkan Cirik, Eduard Hovy, and Louis-Philippe Morency, “Visualizing and Understanding Cur-  
467 riculum Learning for Long Short-Term Memory Networks,” *arXiv preprint arXiv:1611.06204*,  
468 2016.
- 469 [45] Radu Tudor Ionescu, Bogdan Alexe, Marius Leordeanu, Marius Popescu, Dim P. Papadopoulos,  
470 and Vittorio Ferrari, “How Hard Can It Be? Estimating the Difficulty of Visual Search in an  
471 Image,” in *Proceedings of CVPR*, 2016, pp. 2157–2166.
- 472 [46] Chen Gong, Dacheng Tao, Stephen J. Maybank, Wei Liu, Guoliang Kang, and Jie Yang, “Multi-  
473 Modal Curriculum Learning for Semi-Supervised Image Classification,” *IEEE Transactions on  
474 Image Processing*, vol. 25, no. 7, pp. 3249–3260, 2016.
- 475 [47] Guy Hacohen and Daphna Weinshall, “On The Power of Curriculum Learning in Training Deep  
476 Networks,” in *Proceedings of ICML*, 2019, pp. 2535–2544.
- 477 [48] M. Kumar, Benjamin Packer, and Daphne Koller, “Self-Paced Learning for Latent Variable  
478 Models,” in *Proceedings of NIPS*, 2010, vol. 23, pp. 1189–1197.
- 479 [49] Maoguo Gong, Hao Li, Deyu Meng, Qiguang Miao, and Jia Liu, “Decomposition-based evolu-  
480 tionary multiobjective optimization to self-paced learning,” *IEEE Transactions on Evolutionary  
481 Computation*, vol. 23, no. 2, pp. 288–302, 2019.
- 482 [50] Yanbo Fan, Ran He, Jian Liang, and Bao-Gang Hu, “Self-Paced Learning: An Implicit  
483 Regularization Perspective,” in *Proceedings of AAAI*, 2017, pp. 1877–1883.

- 484 [51] Hao Li, Maoguo Gong, Deyu Meng, and Qiguang Miao, “Multi-Objective Self-Paced Learning,”  
485 in *Proceedings of AAAI*, 2016, pp. 1802–1808.
- 486 [52] Sanping Zhou, Jinjun Wang, Deyu Meng, Xiaomeng Xin, Yubing Li, Yihong Gong, and Nanning  
487 Zheng, “Deep self-paced learning for person re-identification,” *Pattern Recognition*, vol. 76, pp.  
488 739–751, 2018.
- 489 [53] Lu Jiang, Deyu Meng, Qian Zhao, Shiguang Shan, and Alexander G. Hauptmann, “Self-Paced  
490 Curriculum Learning,” in *Proceedings of AAAI*, 2015, pp. 2694–2700.
- 491 [54] Nicolae-Catalin Ristea and Radu Tudor Ionescu, “Self-paced ensemble learning for speech and  
492 audio classification,” in *Proceedings of INTERSPEECH*, 2021, pp. 2836–2840.
- 493 [55] Fan Ma, Deyu Meng, Qi Xie, Zina Li, and Xuanyi Dong, “Self-paced co-training,” in  
494 *Proceedings of ICML*, 2017, vol. 70, pp. 2275–2284.
- 495 [56] Lu Jiang, Deyu Meng, Shou-I Yu, Zhenzhong Lan, Shiguang Shan, and Alexander G. Haupt-  
496 mann, “Self-Paced Learning with Diversity,” in *Proceedings of NIPS*, 2014, pp. 2078–2086.
- 497 [57] Min Zhang, Zhongwei Yu, Hai Wang, Hongbo Qin, Wei Zhao, and Yan Liu, “Automatic Digital  
498 Modulation Classification Based on Curriculum Learning,” *Applied Sciences*, vol. 9, no. 10,  
499 2019.
- 500 [58] Lijun Wu, Fei Tian, Yingce Xia, Yang Fan, Tao Qin, Lai Jian-Huang, and Tie-Yan Liu, “Learning  
501 to Teach with Dynamic Loss Functions,” in *Proceedings of NeurIPS*, 2018, vol. 31, pp. 6467–  
502 6478.
- 503 [59] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,  
504 Omer Levy, and Samuel Bowman, “SuperGLUE: A Stickier Benchmark for General-Purpose  
505 Language Understanding Systems,” in *Proceedings of NeurIPS*, 2019, vol. 32, pp. 3266–3280.
- 506 [60] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang, “SQuAD: 100,000+  
507 Questions for Machine Comprehension of Text,” in *Proceedings of EMNLP*, 2016, pp. 2383–  
508 2392.
- 509 [61] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk,  
510 and Quoc V Le, “SpecAugment: A simple data augmentation method for automatic speech  
511 recognition,” *Proceedings of INTERSPEECH*, pp. 2613–2617, 2019.

## 512 Checklist

513 The checklist follows the references. Please read the checklist guidelines carefully for information on  
514 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
515 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
516 the appropriate section of your paper or providing a brief inline description. For example:

- 517 • Did you include the license to the code and datasets? **[Yes]**
- 518 • Did you include the license to the code and datasets? **[No]** The code and the data are  
519 proprietary.
- 520 • Did you include the license to the code and datasets? **[N/A]**

521 Please do not modify the questions and only use the provided macros for your answers. Note that the  
522 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
523 block and only keep the Checklist section heading above along with the questions/answers below.

524 1. For all authors...

- 525 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
526 contributions and scope? **[Yes]**
- 527 (b) Did you describe the limitations of your work? **[No]** We did not identify any significant  
528 limitations of LeRaC.

- 529 (c) Did you discuss any potential negative societal impacts of your work? [No] We did not  
530 identify any negative societal impacts.
- 531 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
532 them? [Yes]
- 533 2. If you are including theoretical results...
- 534 (a) Did you state the full set of assumptions of all theoretical results? [N/A]  
535 (b) Did you include complete proofs of all theoretical results? [N/A]
- 536 3. If you ran experiments...
- 537 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
538 perimental results (either in the supplemental material or as a URL)? [Yes] In the  
539 supplementary.
- 540 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
541 were chosen)? [Yes] See Table 1.
- 542 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
543 ments multiple times)? [Yes] We report average accuracy rates  $\pm$  standard deviations  
544 over 5 runs.
- 545 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
546 of GPUs, internal cluster, or cloud provider)? [Yes] For some cases, please see Figure  
547 1.
- 548 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 549 (a) If your work uses existing assets, did you cite the creators? [Yes]  
550 (b) Did you mention the license of the assets? [No] But we provided the download link of  
551 the used code.
- 552 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
553 We will release our code after acceptance.
- 554 (d) Did you discuss whether and how consent was obtained from people whose data you're  
555 using/curating? [N/A]
- 556 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
557 information or offensive content? [N/A]
- 558 5. If you used crowdsourcing or conducted research with human subjects...
- 559 (a) Did you include the full text of instructions given to participants and screenshots, if  
560 applicable? [N/A]
- 561 (b) Did you describe any potential participant risks, with links to Institutional Review  
562 Board (IRB) approvals, if applicable? [N/A]
- 563 (c) Did you include the estimated hourly wage paid to participants and the total amount  
564 spent on participant compensation? [N/A]