



清华大学  
Tsinghua University



**ICML**  
International Conference  
On Machine Learning

# CurBench: Curriculum Learning Benchmark

Yuwei Zhou, Zirui Pan, Xin Wang, Hong Chen, Haoyang Li, Yanwen Huang,

Zhixiao Xiong, Fangzhou Xiong, Peiyang Xu, Shengnan Liu, Wenwu Zhu

Media and Network Lab, Tsinghua University

# Background

## Curriculum Learning

- Curriculum learning is a training paradigm where machine learning models are trained **in a meaningful order**, inspired by the way humans learn curricula.
- It brings the advantage of **enhancing model generalization** and **accelerating convergence speed**.

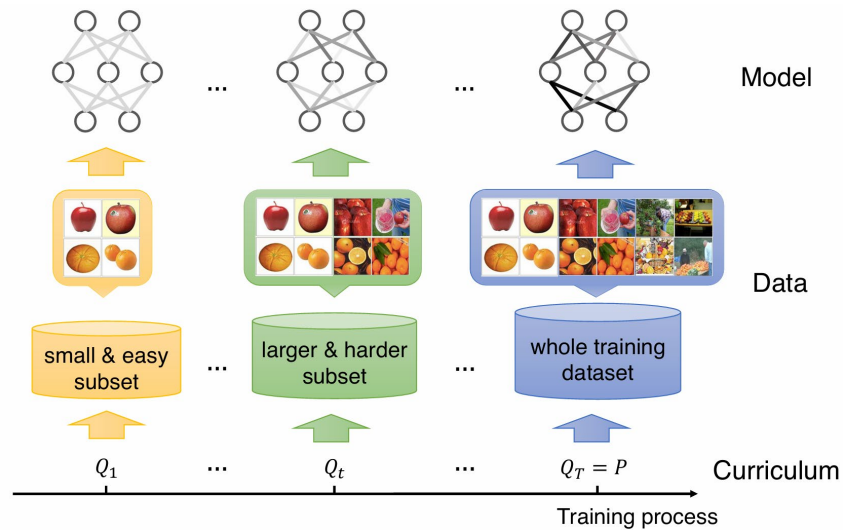


Illustration of Curriculum Learning Concept from [1].

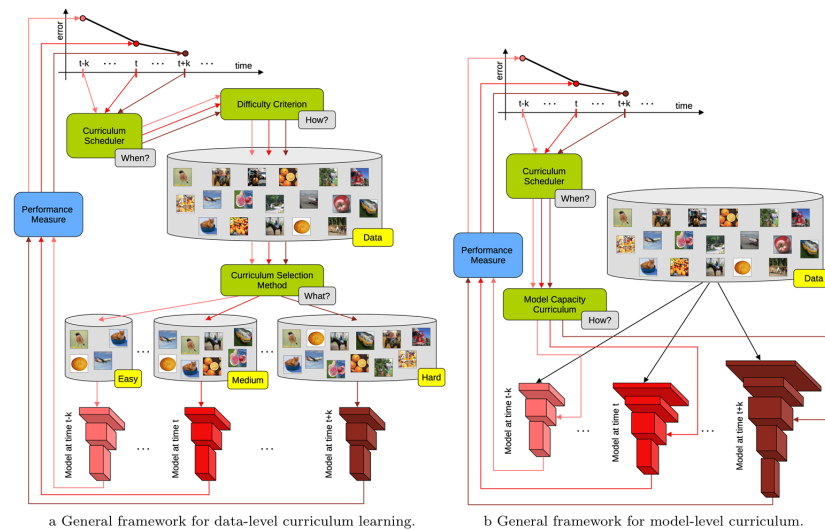


Illustration of Curriculum Learning Framework from [2].

[1] A Survey on Curriculum Learning. TPAMI 2021.

[2] Curriculum Learning: A Survey. IJCV 2022.

# Problem

## No Benchmark for Curriculum Learning

- As new curriculum learning methods continue to emerge, it remains an open issue to benchmark them.
- The increasing number of works pose challenges in terms of comparison and evaluation, mainly due to **the differences in the experimental setups** including datasets, backbone models, and settings.

---

### Methods of Comparison

### The Same

### The Difference

DCL v.s. DDS

WideResNet-28-10

CIFAR-100 v.s. CIFAR-10

DIHCL v.s. CBS

ImageNet

ResNet-50 v.s. ResNet-18

MCL v.s. LRE

MNIST and LeNet

Standard v.s. Imbalance

---

# Related Works

## Summative Work on Curriculum Learning

- **From a theoretical perspective:**

- General Curriculum Learning:

- A Survey on Curriculum Learning. TPAMI 2021.
- Curriculum Learning: A Survey. IJCV 2022.

- Curriculum Learning for Reinforcement Learning:

- Curriculum Learning for Reinforcement Learning Domains: A Framework and Survey. JMLR 2020.
- Automatic Curriculum Learning For Deep RL: A Short Survey. IJCAI 2020.

- Curriculum Learning for Graph Machine Learning:

- Curriculum graph machine learning: A survey. IJCAI 2023.

- **From an empirical perspective:**

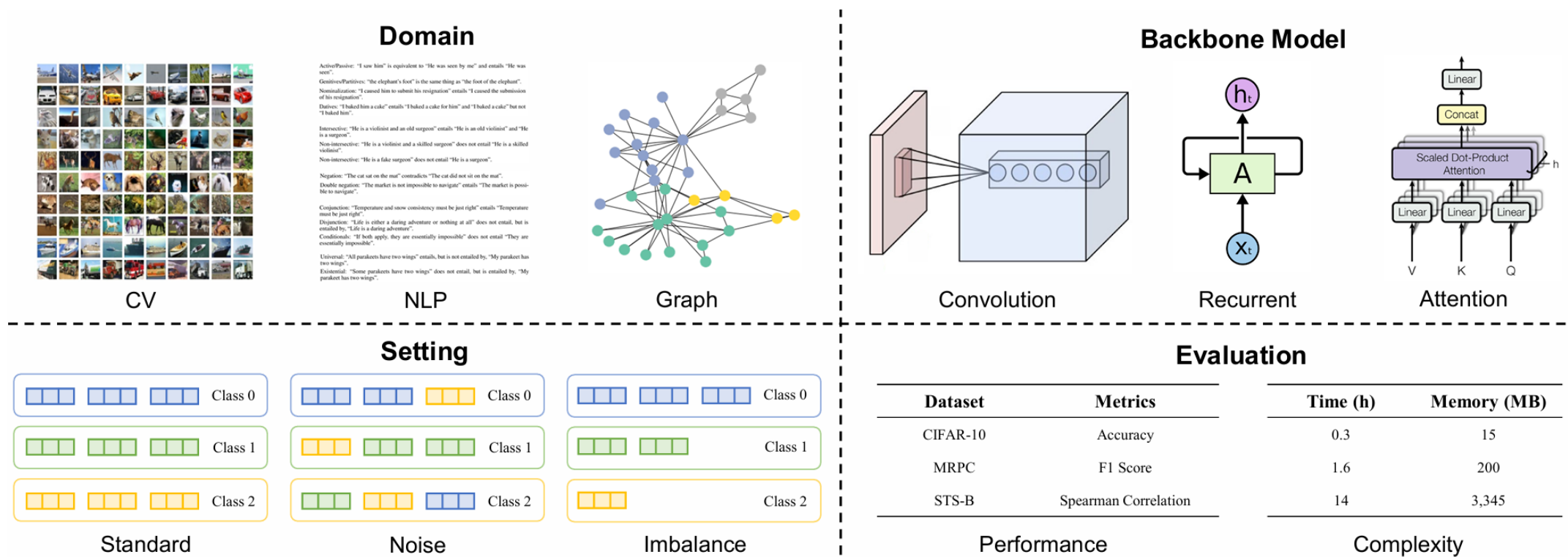
- Curriculum Learning Library:

- CurML: A Curriculum Machine Learning Library. ACMML 2022.

# CurBench

## Outline

- CurBench includes 15 datasets spanning 3 research domains, 9 backbone models, 3 training settings, and 2 evaluation dimensions, with a toolkit for reproducing 15 core curriculum learning methods.



## Dataset

Domain	Dataset	Setting	Training	Validation	Test	Class	Metrics
CV	CIFAR-10	Standard / Noise-0.4	45,000	5,000	10,000	10	Accuracy
		Imbalance-50	12,536	5,000	10,000	10	Accuracy
	CIFAR-100	Standard / Noise-0.4	45,000	5,000	10,000	100	Accuracy
		Imbalance-50	12,536	5,000	10,000	100	Accuracy
	Tiny-ImageNet	Standard / Noise-0.4	90,000	10,000	10,000	200	Accuracy
		Imbalance-50	22,700	10,000	10,000	200	Accuracy
NLP	RTE	Standard / Noise-0.4	2,490	277	-	2	Accuracy
	MRPC	Standard / Noise-0.4	3,668	408	-	2	F1 Score
	STS-B	Standard / Noise-0.4	5,749	1,500	-	6	Spearman
	CoLA	Standard / Noise-0.4	8,551	1,043	-	2	Matthews
	SST-2	Standard / Noise-0.4	67,349	872	-	2	Accuracy
	QNLI	Standard / Noise-0.4	104,743	5,463	-	2	Accuracy
	QQP	Standard / Noise-0.4	363,846	40,430	-	2	F1 Score
	MNLI-(m/mm)	Standard / Noise-0.4	392,702	9,815/9,832	-	3	Accuracy
Graph	MUTAG	Standard / Noise-0.4	150	19	19	2	Accuracy
	PROTEINS	Standard / Noise-0.4	890	111	112	2	Accuracy
	NCI1	Standard / Noise-0.4	3,288	411	411	2	Accuracy
	ogbg-molhiv	Standard / Noise-0.4	32,901	4,113	4,113	2	ROC-AUC

# CurBench

## Model

Domain	Model	Mechanism	Parameters
CV	LeNet	Convolution	$\sim 0.07\text{M}$
	ResNet-18	Convolution	$\sim 11.2\text{M}$
	ViT	Attention	$\sim 9.6\text{M}$
NLP	LSTM	Recurrent	$\sim 10.4\text{M}$
	BERT	Attention	$\sim 109\text{M}$
	GPT2	Attention	$\sim 124\text{M}$
Graph	GCN	Convolution	$\sim 0.01\text{M}$
	GAT	Attention	$\sim 0.14\text{M}$
	GIN	Isomorphism	$\sim 0.01\text{M}$

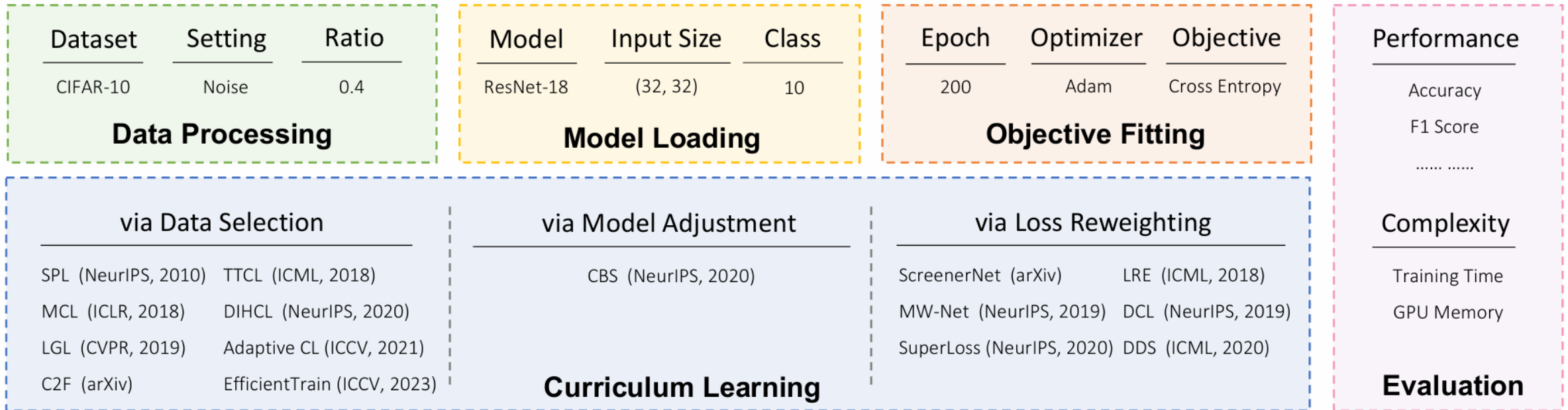
## Setting

- Standard: No additional data processing.
- Noise- $p$ :  $p\%$  data samples are independently attached with random incorrect labels.
- Imbalance- $r$ : A ratio of  $r$  between the number of samples in the largest class and that in the smallest class in a long-tailed dataset where the number of samples for each class follows a geometric sequence.

## Evaluation

- Performance: We report the average and standard deviation of the metric over 5 runs.
- Complexity: We record the training time and maximum memory consumption on the same GPU device.

## Toolkit





# Experiment

## Main Results on CV and Graph Datasets

- There has been no such method that outperforms others all the time, and the effectiveness depends on specific scenarios.

	CIFAR-10			CIFAR-100			Tiny-ImageNet		
	Standard	Noise-0.4	Imbalance-50	Standard	Noise-0.4	Imbalance-50	Standard	Noise-0.4	Imbalance-50
LeNet	69.95 <sub>1.00</sub>	65.02 <sub>1.12</sub>	44.93 <sub>0.56</sub>	35.46 <sub>0.70</sub>	29.59 <sub>0.40</sub>	19.57 <sub>0.64</sub>	22.08 <sub>0.61</sub>	18.63 <sub>0.43</sub>	11.65 <sub>0.30</sub>
LeNet + CL	<b>70.43</b> <sub>0.41</sub>	<b>65.93</b> <sub>0.57</sub>	<b>45.28</b> <sub>0.56</sub>	<b>35.63</b> <sub>0.78</sub>	<b>30.87</b> <sub>0.48</sub>	<b>19.74</b> <sub>0.17</sub>	<b>22.83</b> <sub>0.44</sub>	<b>19.91</b> <sub>0.26</sub>	<b>12.36</b> <sub>0.47</sub>
ResNet-18	92.33 <sub>0.16</sub>	82.75 <sub>2.06</sub>	75.49 <sub>0.87</sub>	69.97 <sub>0.27</sub>	52.14 <sub>0.39</sub>	42.57 <sub>0.68</sub>	51.41 <sub>1.74</sub>	39.42 <sub>0.21</sub>	28.83 <sub>0.38</sub>
ResNet-18 + CL	<b>92.88</b> <sub>0.23</sub>	<b>86.92</b> <sub>0.20</sub>	<b>76.43</b> <sub>0.96</sub>	<b>71.31</b> <sub>0.14</sub>	<b>58.56</b> <sub>0.60</sub>	<b>43.47</b> <sub>0.43</sub>	<b>53.61</b> <sub>0.48</sub>	<b>43.64</b> <sub>0.72</sub>	<b>30.82</b> <sub>0.36</sub>
ViT	79.90 <sub>0.38</sub>	64.19 <sub>0.51</sub>	52.12 <sub>0.81</sub>	51.05 <sub>0.62</sub>	35.25 <sub>0.24</sub>	26.05 <sub>0.52</sub>	38.16 <sub>0.53</sub>	24.90 <sub>0.26</sub>	17.15 <sub>0.31</sub>
ViT + CL	<b>80.66</b> <sub>0.27</sub>	<b>69.83</b> <sub>0.53</sub>	<b>52.85</b> <sub>0.81</sub>	<b>51.93</b> <sub>0.64</sub>	<b>39.15</b> <sub>0.30</sub>	<b>26.40</b> <sub>0.34</sub>	<b>38.92</b> <sub>0.53</sub>	<b>29.76</b> <sub>0.34</sub>	<b>17.47</b> <sub>0.14</sub>

	MUTAG		PROTEINS		NCI1		ogbg-molhiv	
	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4
GCN	73.68 <sub>2.11</sub>	66.31 <sub>7.14</sub>	70.71 <sub>4.20</sub>	63.57 <sub>6.45</sub>	69.59 <sub>1.23</sub>	55.23 <sub>3.21</sub>	75.84 <sub>1.02</sub>	64.29 <sub>4.55</sub>
GCN + CL	<b>74.74</b> <sub>3.94</sub>	<b>71.58</b> <sub>5.37</sub>	<b>73.21</b> <sub>4.41</sub>	<b>71.61</b> <sub>6.62</sub>	<b>71.39</b> <sub>1.29</sub>	<b>67.98</b> <sub>2.01</sub>	<b>77.41</b> <sub>1.15</sub>	<b>72.81</b> <sub>1.14</sub>
GAT	69.47 <sub>6.14</sub>	65.26 <sub>5.37</sub>	64.46 <sub>2.96</sub>	65.71 <sub>9.13</sub>	56.74 <sub>2.86</sub>	53.77 <sub>2.12</sub>	68.07 <sub>2.34</sub>	65.37 <sub>2.66</sub>
GAT + CL	<b>72.63</b> <sub>8.42</sub>	<b>69.47</b> <sub>10.21</sub>	<b>69.82</b> <sub>7.13</sub>	<b>69.11</b> <sub>3.77</sub>	<b>59.37</b> <sub>1.59</sub>	<b>55.67</b> <sub>4.70</sub>	<b>72.64</b> <sub>1.16</sub>	<b>66.73</b> <sub>1.84</sub>
GIN	86.84 <sub>7.90</sub>	78.95 <sub>3.72</sub>	74.11 <sub>4.24</sub>	69.82 <sub>1.73</sub>	79.32 <sub>1.40</sub>	60.24 <sub>3.92</sub>	74.72 <sub>1.36</sub>	63.07 <sub>3.73</sub>
GIN + CL	<b>88.42</b> <sub>2.10</sub>	<b>81.58</b> <sub>4.56</sub>	<b>77.14</b> <sub>4.88</sub>	<b>73.93</b> <sub>1.82</sub>	<b>82.04</b> <sub>1.90</sub>	<b>62.14</b> <sub>6.47</sub>	<b>76.53</b> <sub>1.97</sub>	<b>65.53</b> <sub>1.61</sub>

# Experiment

## Main Results on NLP Datasets

- There has been no such method that outperforms others all the time, and the effectiveness depends on specific scenarios.

	RTE		MRPC		STS-B		CoLA	
	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4
LSTM	52.95 <sub>1.34</sub>	53.43 <sub>1.77</sub>	81.43 <sub>0.14</sub>	81.22 <sub>0.00</sub>	12.73 <sub>0.72</sub>	10.90 <sub>1.19</sub>	11.29 <sub>1.27</sub>	3.27 <sub>1.68</sub>
LSTM + CL	<b>53.07</b> <sub>1.29</sub>	<b>54.22</b> <sub>1.77</sub>	<b>81.54</b> <sub>0.18</sub>	<b>81.24</b> <sub>0.05</sub>	<b>14.11</b> <sub>2.21</sub>	<b>11.75</b> <sub>1.61</sub>	<b>12.65</b> <sub>1.21</sub>	<b>8.55</b> <sub>2.10</sub>
BERT	64.62 <sub>3.33</sub>	54.22 <sub>3.14</sub>	88.54 <sub>0.45</sub>	81.89 <sub>0.83</sub>	85.26 <sub>0.22</sub>	80.71 <sub>1.01</sub>	57.39 <sub>1.30</sub>	32.35 <sub>0.79</sub>
BERT + CL	<b>66.35</b> <sub>1.76</sub>	<b>56.32</b> <sub>5.04</sub>	<b>88.69</b> <sub>1.24</sub>	<b>81.94</b> <sub>0.55</sub>	<b>85.42</b> <sub>0.22</sub>	<b>81.31</b> <sub>0.25</sub>	<b>57.80</b> <sub>1.96</sub>	<b>45.79</b> <sub>1.64</sub>
GPT2	65.34 <sub>1.95</sub>	52.92 <sub>4.49</sub>	85.49 <sub>0.86</sub>	78.23 <sub>1.72</sub>	76.44 <sub>1.20</sub>	69.65 <sub>1.85</sub>	37.00 <sub>3.72</sub>	5.86 <sub>1.69</sub>
GPT2 + CL	<b>66.35</b> <sub>2.10</sub>	<b>57.40</b> <sub>3.39</sub>	<b>86.29</b> <sub>0.36</sub>	<b>82.55</b> <sub>0.88</sub>	<b>80.82</b> <sub>1.39</sub>	<b>71.57</b> <sub>1.74</sub>	<b>39.95</b> <sub>3.16</sub>	<b>12.54</b> <sub>2.75</sub>

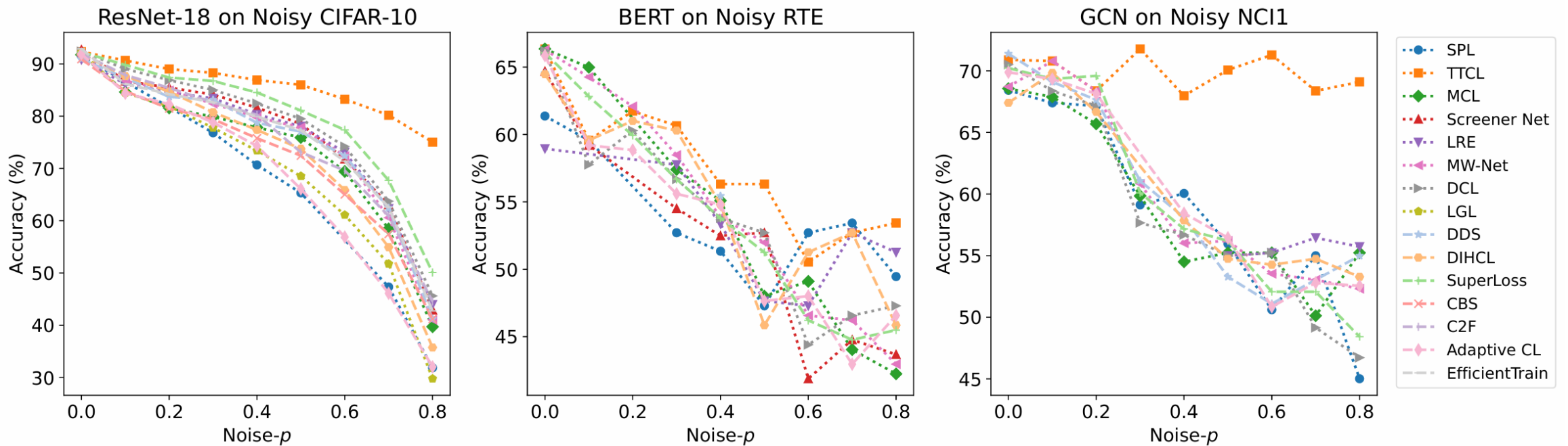
  

	SST-2		QNLI		QQP		MNLI-(m/mm)	
	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4	Standard	Noise-0.4
LSTM	81.67 <sub>0.85</sub>	64.36 <sub>1.12</sub>	50.54 <sub>0.00</sub>	50.62 <sub>0.16</sub>	75.69 <sub>0.27</sub>	60.72 <sub>0.79</sub>	61.38 <sub>0.30</sub> / 61.21 <sub>0.45</sub>	44.41 <sub>0.51</sub> / 44.83 <sub>0.90</sub>
LSTM + CL	<b>82.87</b> <sub>0.88</sub>	<b>78.58</b> <sub>1.64</sub>	<b>51.02</b> <sub>0.46</sub>	<b>50.83</b> <sub>0.45</sub>	<b>75.73</b> <sub>0.21</sub>	<b>66.47</b> <sub>0.72</sub>	<b>62.47</b> <sub>0.36</sub> / <b>62.33</b> <sub>0.42</sub>	<b>58.59</b> <sub>0.54</sub> / <b>58.50</b> <sub>0.64</sub>
BERT	92.66 <sub>0.28</sub>	87.22 <sub>0.82</sub>	91.21 <sub>0.24</sub>	81.21 <sub>0.76</sub>	88.05 <sub>0.12</sub>	76.23 <sub>0.48</sub>	83.89 <sub>0.31</sub> / 84.38 <sub>0.29</sub>	78.65 <sub>0.70</sub> / 79.21 <sub>0.62</sub>
BERT + CL	<b>92.82</b> <sub>0.16</sub>	<b>91.25</b> <sub>0.59</sub>	<b>91.49</b> <sub>0.13</sub>	<b>89.45</b> <sub>0.44</sub>	<b>88.16</b> <sub>0.13</sub>	<b>84.50</b> <sub>0.25</sub>	<b>84.27</b> <sub>0.07</sub> / <b>84.40</b> <sub>0.42</sub>	<b>81.73</b> <sub>0.31</sub> / <b>82.25</b> <sub>0.40</sub>
GPT2	91.95 <sub>0.49</sub>	85.83 <sub>0.57</sub>	87.92 <sub>0.31</sub>	78.72 <sub>0.37</sub>	86.00 <sub>0.23</sub>	75.40 <sub>0.84</sub>	81.53 <sub>0.21</sub> / 82.40 <sub>0.21</sub>	76.56 <sub>0.15</sub> / 77.69 <sub>0.15</sub>
GPT2 + CL	<b>92.25</b> <sub>0.42</sub>	<b>90.34</b> <sub>0.53</sub>	<b>88.17</b> <sub>0.67</sub>	<b>84.00</b> <sub>0.70</sub>	<b>86.68</b> <sub>0.16</sub>	<b>82.16</b> <sub>0.35</sub>	<b>81.90</b> <sub>0.23</sub> / <b>82.59</b> <sub>0.35</sub>	<b>78.36</b> <sub>0.19</sub> / <b>79.62</b> <sub>0.44</sub>

# Experiment

## Results in Noise Settings

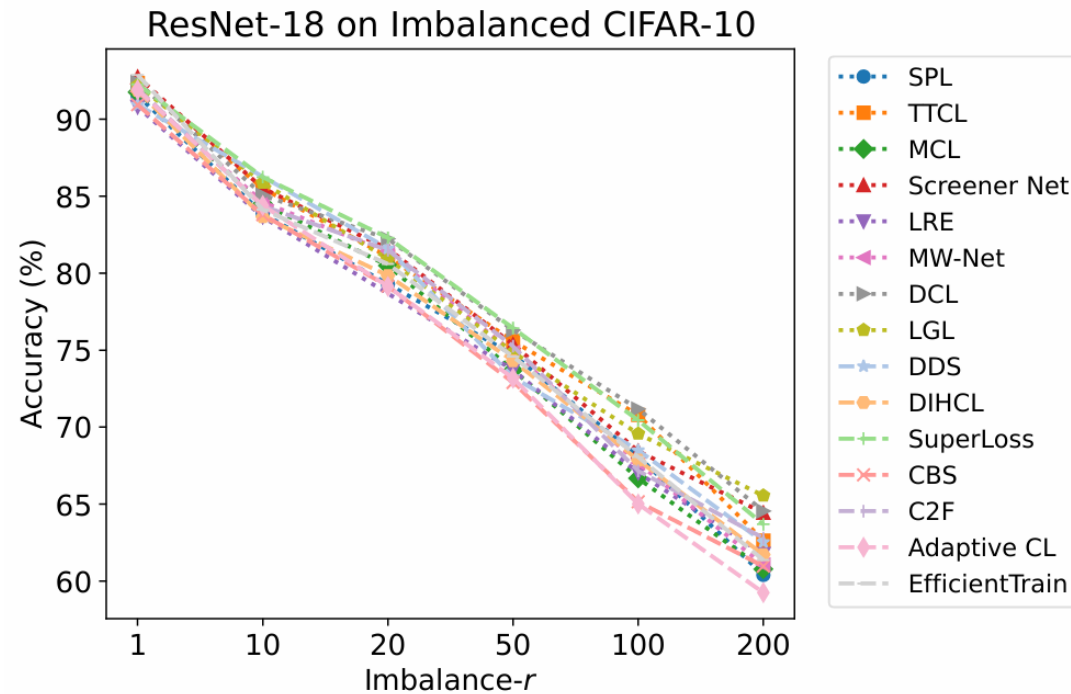
- Methods by teacher transferring have edges in noise settings.



# Experiment

## Results in Imbalance Settings

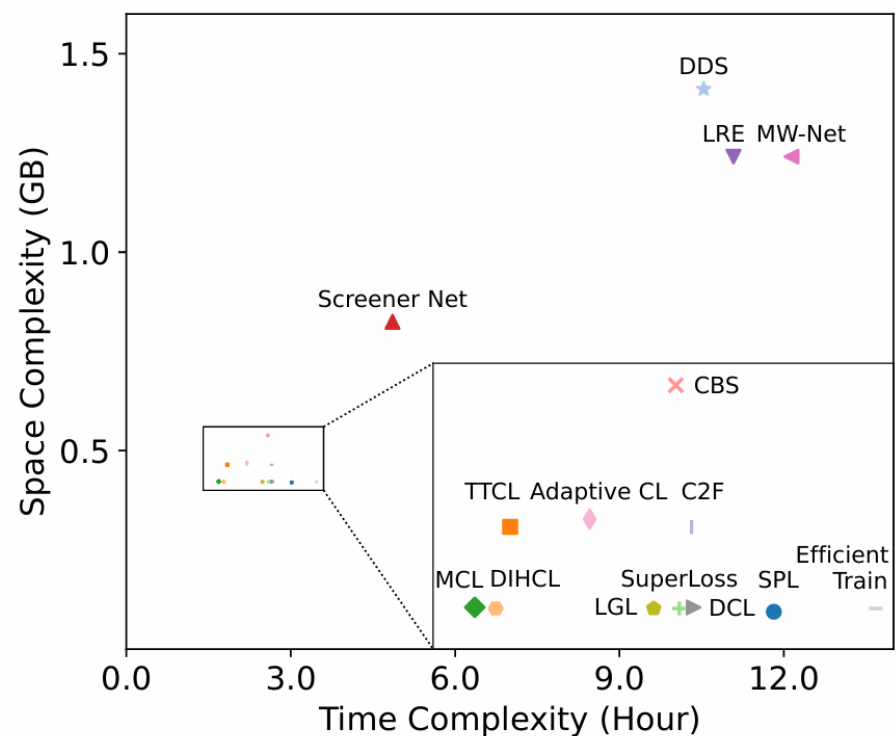
- All methods achieve similar performances under different imbalance ratios.
- Methods by reweighting perform relatively well in imbalance settings.



# Experiment

## Time and Space Complexity

- Methods involving gradient calculation and extra learnable networks generally have higher time and space complexity.



# Summary

## Findings

- 1) There has been no such method that outperforms others all the time, and the effectiveness depends on specific scenarios.
- 2) Curriculum learning brings more significant improvements in noise settings than in standard and imbalance ones.
- 3) Methods by teacher transferring have edges in noise settings, while methods by reweighting perform relatively well in imbalance settings.
- 4) Methods involving gradient calculation and extra learnable networks generally have higher time and space complexity.

# Summary

## Contributions

- 1) We propose CurBench, the first benchmark on curriculum learning to the best of our knowledge.
- 2) We conduct extensive experiments to impartially evaluate and compare the performance and complexity of existing curriculum learning methods under various experimental setups.
- 3) We make in-depth analyses and demonstrate intriguing observations on curriculum learning based on empirical results derived from CurBench.



清华大学  
Tsinghua University



**ICML**  
International Conference  
On Machine Learning

# CurBench: Curriculum Learning Benchmark

Yuwei Zhou, Zirui Pan, Xin Wang, Hong Chen, Haoyang Li, Yanwen Huang,

Zhixiao Xiong, Fangzhou Xiong, Peiyang Xu, Shengnan Liu, Wenwu Zhu

Media and Network Lab, Tsinghua University