

# STAR: Stability-Inducing Weight Perturbation for Continual Learning

Masih Eskandar, Tooba Imtiaz, Davin Hill, Zifeng Wang, Jennifer Dy

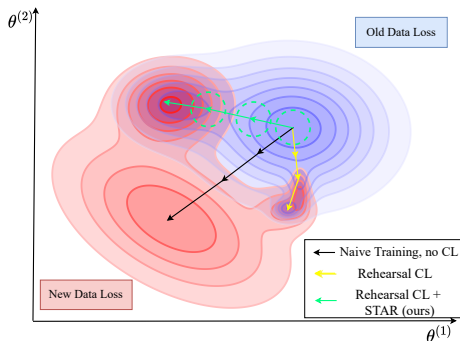
Northeastern University

# Motivation

- Humans learn sequentially without forgetting.
- Neural networks suffer from **catastrophic forgetting**.
- Rehearsal-based methods help but are limited by buffer size.
- Goal: Improve stability of model predictions over time.

# Key Idea: STAR

- Introduce **STAR**: Stability-Inducing Parameter-space Regularization.
- Promotes consistency of model outputs under **worst-case local parameter perturbations**.
- Plug-and-play with any rehearsal-based continual learning method.



# How STAR Works

- 1 Identify correctly classified buffered samples.
- 2 Apply perturbation  $\delta$  that maximizes KL divergence of outputs.
- 3 Use this worst-case  $\delta$  to regularize training:

$$\mathcal{L}_{\text{STAR}} = \max_{\|\delta\| \leq d} \sum_{(x,y) \in M^*} \text{KL}(q_{\theta}(x) \parallel q_{\theta+\delta}(x))$$

- 4 Final objective:

$$\mathcal{L} = \mathcal{L}_{\text{CL}} + \lambda \mathcal{L}_{\text{STAR}}$$

# Results Summary

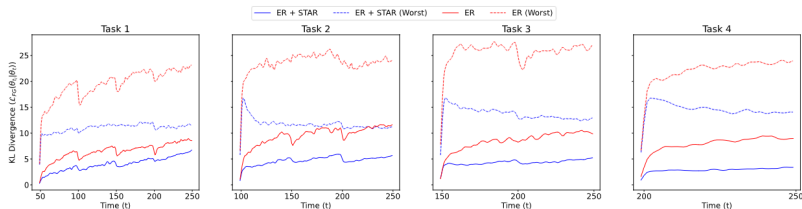
- STAR improves accuracy by up to **15%** across methods and datasets.
- Effective across buffer sizes and baselines (ER, DER++, ER-ACE, X-DER).
- Especially beneficial when buffer is small.

Table 1: Comparison of adding STAR to baseline rehearsal-based CL methods in terms of average accuracy.

Method	Split-Cifar10			Split-Cifar100			Split-miniImageNet		
Sequential	19.67			9.29			4.51		
Joint	92.38			73.29			53.55		
Buffer Size	100	200	500	200	500	2000	1000	2000	5000
ER	36.39	44.79	57.74	14.35	19.66	36.76	8.37	16.49	24.17
+ STAR (ours)	<b>51.5</b>	<b>59.3</b>	<b>70.70</b>	<b>19.64</b>	<b>29.64</b>	<b>44.65</b>	<b>11.83</b>	<b>16.64</b>	<b>25.83</b>
ER-ACE	52.95	61.25	71.16	29.22	38.01	49.95	17.95	22.60	27.92
+ STAR (ours)	<b>60.69</b>	<b>67.58</b>	<b>75.44</b>	<b>30.38</b>	<b>40.20</b>	<b>51.67</b>	<b>21.06</b>	<b>24.9</b>	<b>31.01</b>
DER++	57.65	64.88	72.70	25.11	37.13	52.08	18.02	23.44	30.43
+ STAR (ours)	<b>61.76</b>	<b>68.60</b>	<b>76.52</b>	<b>27.64</b>	<b>39.77</b>	<b>53.24</b>	<b>22.4</b>	<b>28.19</b>	<b>33.36</b>
X-DER (RPC)	52.75	58.48	64.77	37.23	<b>48.53</b>	57.00	23.19	26.38	29.91
+ STAR (ours)	<b>58.85</b>	<b>65.94</b>	<b>69.19</b>	<b>38.15</b>	47.56	<b>57.55</b>	<b>24.6</b>	<b>27.95</b>	<b>32.6</b>

# Ablations and Insights

- Gradient-based perturbations outperform random ones.
- Using only buffer samples in the STAR loss is best.
- Empirically validate that STAR reduces local-worst case change in distribution



# Conclusion

- STAR is a simple yet powerful regularization strategy for CL.
- Enhances rehearsal-based methods without needing task boundaries.
- Open-sourced: [github.com/Gnomy17/STAR\\_CL](https://github.com/Gnomy17/STAR_CL)