

#### Motivations

#### How to design a 'fast' AT with improved stability mitigated catastrophic overfitting, and theoretical guarantees?

## Fast Robust Training: Not Enough!

- Low Stability;
- Robust Catastrophic Overfitting;
- Robustness at Cost of Sharp Drop in Accuracy.

### Bi-level Optimization: A Fresh Perspectiv

Standard min-max formulation for adversarial training

$$\min_{\boldsymbol{\theta}} \mathbf{E}_{(\boldsymbol{x},t)\sim D} \left[ \max_{|\boldsymbol{\delta}|_{\infty} \leq \epsilon} \ell_{\mathrm{tr}}(\boldsymbol{\theta}; \boldsymbol{x} + \boldsymbol{\delta}, t) \right]$$

Bi-level Optimization (BLO) framework:  $\min_{\boldsymbol{\theta}} \boldsymbol{\ell}_{tr}(\boldsymbol{\theta}, \boldsymbol{\delta}^{*}(\boldsymbol{\theta}))$ 

s.t.  $\delta^*(\theta) = \operatorname{argmin}_{\delta \in C} \ell_{atk}(\theta, \delta)$ 

- ✓ Customizable lower-level objectives:  $\ell_{atk} \neq -\ell_{tr}$ !
- $\checkmark$  BLO is a generalization form of min-max optimiza
- Presence of implicit gradient (IG): the 'fingerprint'  $\checkmark$ The upper-level gradient calculation:

$$\frac{\mathrm{d}\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\mathrm{d}\boldsymbol{\theta}} = \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\partial\boldsymbol{\theta}} + \frac{\frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}}}{\underbrace{\mathrm{d}\boldsymbol{\theta}}_{\mathrm{IG}}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\partial\boldsymbol{\theta}} + \frac{\frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}}}{\underbrace{\mathrm{d}\boldsymbol{\theta}}_{\mathrm{IG}}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\partial\boldsymbol{\theta}} + \frac{\frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}}}{\underbrace{\mathrm{d}\boldsymbol{\theta}}_{\mathrm{IG}}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\partial\boldsymbol{\theta}} + \frac{\frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}}}{\underbrace{\mathrm{d}\boldsymbol{\theta}}_{\mathrm{IG}}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta}))}{\partial\boldsymbol{\theta}} + \frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta})}{\partial\boldsymbol{\theta}} + \frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}} \frac{\partial\ell_{\mathrm{tr}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*}(\boldsymbol{\theta})}{\partial\boldsymbol{\theta}} + \frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}} + \frac{\mathrm{d}\boldsymbol{\delta}^{*}$$

**IG**:  $\delta^*(\theta)$  is an implicit function of  $\theta$ 

BLO is hard to solve, lower-level constraint makes

[1] Eric Wong et al. "Fast is Better Than Free: Revisiting Adv [2] Maksym Andriushchenko et al. "Understanding and Impro

# **Revisiting and Advancing Fast Adversarial Training through the Lens of Bi-Level Optimization**

Yihua Zhang<sup>1,\*</sup>, Guanhua Zhang<sup>2,\*</sup>, Prashant Khanduri<sup>3</sup>, Mingyi Hong<sup>3</sup>, Shiyu Chang<sup>2</sup>, Sijia Liu<sup>1,4</sup> <sup>1</sup>Michigan State University, <sup>2</sup>UC Santa Barbara, <sup>3</sup>University of Minnesota, <sup>4</sup>MIT-IBM Watson Lab

۲۰	• Customize Your Attack Loss : Fast-BAT!	<b>Table 1.</b> Performance overview of proposed FAST-BAT vs. the baselines FAST-A[1] and FAST-AT-GA [2] on various datasets with PreActResNet-18.								
y,			CIFAR-10. PreActResNet-18							
	$\ell'_{\mathrm{atk}}(\boldsymbol{\theta}, \boldsymbol{\delta}) = \langle \nabla_{\boldsymbol{\delta}=\boldsymbol{z}} \ell_{\mathrm{atk}}(\boldsymbol{\theta}, \boldsymbol{\delta}), \boldsymbol{\delta}-\boldsymbol{z} > + \left(\frac{\pi}{2}\right)   \boldsymbol{\delta}-\boldsymbol{z}  _2^2$	Method	SA (%) ( $\epsilon = 8/255$ )	RA-PGD (%) ( $\epsilon = 8/255$ )	$\begin{vmatrix} RA-AA(\%) \\ \epsilon = 8/255 \end{vmatrix}$	$\begin{vmatrix} SA(\%) \\ (\epsilon = 16/255) \end{vmatrix}$	RA-PGD (%) ( $\epsilon = 16/255$ )	RA-AA (%) ( $\epsilon = 16/255$ )	Time (s/epoch)	
	Fast-BAT Algorithm	FAST-AT FAST-AT-GA	<b>82.39</b> ±0.14 79.71±0.24 79.97±0.12	$45.49 \pm 0.21$ $47.27 \pm 0.22$ $48.83 \pm 0.17$	$\begin{array}{c c} 41.87 \pm 0.15 \\ 43.24 \pm 0.27 \\ 45.19 \pm 0.12 \end{array}$	$\begin{array}{c c} 44.15 \pm 7.27 \\ 58.29 \pm 1.32 \\ 68.16 \pm 0.25 \end{array}$	$21.83 \pm 1.32$ $26.01 \pm 0.16$ $27.69 \pm 0.16$	12.49±0.33 17.97±0.33 <b>18.79</b> ±0.24	23.1 75.3	
	✓ One-step lower-level solver: unique, efficient, closed-form	TASI-DAI				CIFAR-100, PreActResNet-18		10.77±0.24	01.4	
	$\boldsymbol{\delta}^{*}(\boldsymbol{\theta}) = \operatorname{Proj}_{C}(\boldsymbol{z} - (1/\lambda) \nabla_{\boldsymbol{\delta}} \ell_{\mathrm{atk}}(\boldsymbol{\theta}, \boldsymbol{z}))$	FAST-AT FAST-AT-GA <b>FAST-BAT</b>	$\begin{array}{c} \textbf{52.62} \pm 0.18 \\ \textbf{50.06} \pm 0.27 \\ \textbf{50.19} \pm 0.21 \end{array}$	$\begin{array}{c} 24.66{\pm}0.21\\ 24.97{\pm}0.23\\ \textbf{26.49}{\pm}0.20\end{array}$	$\begin{array}{c c} 21.72 \pm 0.17 \\ 21.82 \pm 0.21 \\ \textbf{23.97} \pm 0.15 \end{array}$	$\begin{array}{c c} 21.32{\pm}3.27\\ 32.51{\pm}1.27\\ \textbf{39.29}{\pm}0.53\end{array}$	$\begin{array}{c} 8.62{\pm}1.03\\ 12.27{\pm}0.36\\ \textbf{13.97}{\pm}0.17\end{array}$	$\begin{array}{c} \textbf{6.22}{\pm 0.61} \\ \textbf{9.43}{\pm 0.19} \\ \textbf{11.32}{\pm 0.22} \end{array}$	23.8 77.1 61.6	
	✓ Tractable IG with KKT conditions and Hessian-free assumption.	Tiny-ImageNet, PreActResNet-18								
ve! ng:	$\frac{\mathrm{d}\boldsymbol{\delta}^{*}(\boldsymbol{\theta})^{T}}{\mathrm{d}\boldsymbol{\theta}} = -\left(\frac{1}{2}\right)\nabla_{\boldsymbol{\theta}\boldsymbol{\delta}}\boldsymbol{\ell}_{\mathrm{atk}}(\boldsymbol{\theta},\boldsymbol{\delta}^{*})\boldsymbol{H}_{C}$	FAST-AT FAST-AT-GA <b>FAST-BAT</b>	$\begin{array}{c} 41.37{\pm}3.08\\ 45.52{\pm}0.24\\ \textbf{45.80}{\pm}0.22\end{array}$	$\begin{array}{c} 17.05{\pm}3.25\\ 20.39{\pm}0.19\\ \textbf{21.97}{\pm}0.21\end{array}$	$\begin{array}{c c} 12.31 \pm 2.73 \\ 16.25 \pm 0.17 \\ \textbf{17.64} \pm 0.15 \end{array}$	$\begin{array}{c c} 31.38 \pm 0.19 \\ 29.17 \pm 0.32 \\ \textbf{33.78} \pm 0.23 \end{array}$	$\begin{array}{c} 5.42{\pm}2.17\\ 6.79{\pm}0.27\\ \textbf{8.83}{\pm}0.22\end{array}$	$3.13 \pm 0.24$ $4.27 \pm 0.15$ $5.52 \pm 0.14$	284.6 592.7 572.4	
	Fast-AT[1]: Live by 'Sign'. Die from 'Sign' !	<b>Table 2.</b> Performance comparison on CIFAR-10 with various model architectures.								
	Fast BAT with gradient sign-based linearization:	Model	Meth	od $\epsilon = \frac{1}{\epsilon}$	SA(%) = 8/255)	$RA-PGD(\%)$ $(\epsilon = 8/255)$	SA(%) $(\epsilon = 16/2)$	$\begin{array}{l} \text{RA-P} \\ 55)  (\epsilon = 1 \end{array}$	GD(%) .6/255)	
	$\ell'_{\text{atk}}(\boldsymbol{\theta}, \boldsymbol{\delta}) = < \operatorname{sign}(\nabla_{\boldsymbol{\delta}} \ell_{\text{atk}}(\boldsymbol{\theta}, \boldsymbol{z})), \boldsymbol{\delta} - \boldsymbol{z} > + (\frac{\lambda}{2})   \boldsymbol{\delta} - \boldsymbol{z}  _{2}^{2}$	PARN-50	FAST- FAST-A	AT 73 Г-GA 77	$.15 \pm 6.10$ $.40 \pm 0.81$	$41.03 \pm 2.99$ $46.16 \pm 0.98$	43.86±4. 42.28±6.	<ul> <li>31 22.08</li> <li>69 22.8'</li> <li>22.2</li> </ul>	$22.08 \pm 0.27$ $22.87 \pm 1.25$	
	Fast-AT Algorithm		FAST-E	-A1 <b>83</b> BAT 78	$33.53 \pm 0.17$ $40.17 \pm 0.59$ $78.91 \pm 0.68$ $49.18 \pm 0.35$ $84.39 \pm 0.46$ $45.80 \pm 0.57$ $81.51 \pm 0.38$ $48.29 \pm 0.20$		<b>68.88</b> ±0. <b>69.01</b> ±0.	39       22.3         19 <b>24.5</b>	$22.37 \pm 0.41$ $24.55 \pm 0.06$ $21.99 \pm 0.41$ $23.10 \pm 3.90$	
ation	✓ One-step lower-level solver: <b>unique</b> , <b>efficient</b> , <b>closed-form</b>	WRN-16-8	FAST-	АТ 84 Г-GA 81			49.39±2. 45.95±13	17       21.99         .65       23.19		
' of BL	$O \qquad \delta^*(\boldsymbol{\theta}) = \operatorname{Proj}_{\mathcal{C}}(\boldsymbol{z} - (1/\lambda) \operatorname{sign}(\nabla_{\boldsymbol{\delta}} \ell_{\operatorname{atk}}(\boldsymbol{\theta}, \boldsymbol{z})))$		PGD-2	-AT 85	<b>.52</b> ±0.14	45.47±0.14	<b>72.11</b> ±0.	33 23.6	$1\pm 0.16$	
	✓ IG degrades to zero! BLO degrades to MMO!		FAST-E	3AI   81	<b>.00</b> ±0.54	<b>49.93</b> ±0.36	68.12±0.	47 25.0.	<b>3</b> ±0.44	
$oldsymbol{\delta}^{*}(oldsymbol{ heta}))$	$\frac{\mathrm{d}\boldsymbol{\delta}^*(\boldsymbol{\theta})^T}{\mathbf{d}\boldsymbol{\delta}^*(\boldsymbol{\theta})} = 0$	FAST-AT			FAST-AT-GA			FAST-BAT		
δ	$ \mathbf{d}\boldsymbol{\theta} $ Why not use Fast-AT?			3.00	2.80			2.90		
	$\checkmark$ Sign operation is non-smooth that affects training stability!			2.80		2.	60		2.70	
s it harder! ✓ Sign operation is sub-optimal for linearization!		-2.50-1.25 0.00 1.25	-2.501.250.00	1.25 -2.50 <sub>-1</sub>	.25 0.00 1.25	-2.50 <sup>1.25</sup> 0.00 <sup>1.25</sup>	-2.50-1.25 0.00	1.25 -2.50 <sup>1.25</sup>	0.00 <sup>1.25</sup>	
lversarial Training." ICLR 2020. coving Fast Adversarial Training." NeurIPS 2020.		Figure 1. V	<i>'</i> isualizatio	on of advers	sarial loss l	andscapes (F	ResNet18, C	IFAR10)		

Contact: {zhan1908, líusíjí5}@msu.edu, {guanhua, shíyu}@ucsb.edu, {khand095,mhong}@umn.edu





Code





Figure 2. Robustness against different training attack strengths.



Figure 3. Training curves of Gradient Alignment score. Fast-BAT achieves 'gradient alignment' for free!

