



Securing Artificial Intelligence for Battlefield Effective Robustness (SABER)

LTC Nathaniel D. Bastian, PhD Information Innovation Office (I2O)





8:00 AM - 9:00 AM	Arrival/Check-In	Executive Conference Center
		LTC Nathaniel D. Bastian, PhD
9:00 AM - 9:05 AM	Opening Remarks	Program Manager, DARPA I2O
		Mr. Justin Winburn
9:05 AM - 9:10 AM	Security Review	Program Secuirty Representative
9:10 AM - 9:30 AM	DARPA CMO Review	Contracting Officer Representative
		LTC Nathaniel D. Bastian, PhD
9:30 AM - 10:15 AM	DARPA PM Presentation	Program Manager, DARPA I2O
10:15 AM - 10:30 AM	Break / Attendees submit to Q&A	
		LTC Nathaniel D. Bastian, PhD
10:30 AM - 11:30 AM	Informal Teaming / PM reviews Q&A	Program Manager, DARPA I2O
		LTC Nathaniel D. Bastian, PhD
11:30 AM -12:00 PM	PM Addresses Q&A	Program Manager, DARPA I2O



Operational security risks of AI-enabled battlefield systems are unknown



https://www.defenseone.com/technology/2024/10/new-ai-powered-strike-drone-shows-how-quickly-battlefield-autonomy-evolving/400179/?oref=d1-topic-lander-river

AI-enabled battlefield systems can help improve the speed, quality and accuracy of decisions in the field, providing a decisive advantage

Model Physical Dynamics by Sampling from Distribution $f_{\theta}(x)$ $f_{\theta}(x)$ $f_{\theta}(x)$ $f_{\theta}(x)$ $f_{\theta}(x)$

K. Eykholt et al., "Robust Physical-World attacks on deep learning models," arXiv.org, Jul. 27, 2017. https://arxiv.org/abs/1707.08945

Adversarial AI attacks have not been practically demonstrated in operational settings

Technical Hypothesis:

An AI red team employing SoA PACE techniques and red-teaming tools to execute novel AI kill chain TTPs can operationally assess AI-enabled battlefield systems more effectively than a team without them

Disruption Target:

Build an exemplar AI red team that can continuously integrate and employ emerging counter-AI techniques and tools, establishing a sustainable model for an operational AI red-teaming process

AI = Artificial Intelligence SoA = State-of-the-Art PACE = Physical+Adversarial AI+Cyber+Electronic Warfare TTP = Tactics, Techniques and Procedures

Distribution Statement A (Approved for Public Release, Distribution Unlimited)



AI-enabled battlefield systems are used without knowing the risks

Source: https://www.nytimes.com/2024/07/02/technology/ukraine-war-ai-weapons.htm

Drones



Live video feed taken from an AI-enabled drone, identifying an asset...



... then autonomously delivering its firing effect

Turrets



AI-enabled auto-turrets are giving operators the ability to work remotely with enhanced targeting...



... and can be mounted on autonomous ground robots, relying on auto aim to identify and eliminate targets

Operational AI red-teaming does not exist, making operationally impactful vulnerabilities unknown

DoD = Department of Defense AI = Artificial Intelligence



Adversarial AI access paradigms



Adversarial AI attack categories



No work has coupled physical, adversarial AI, cyber and EW attacks to create AI kill chain TTPs needed for operational AI red-teaming Most previous academic work ignores physically realizable attacks and the full-system pipeline



Example AI-enabled battlefield system pipeline for vulnerability discovery



TTPs = Tactics, Techniques and Procedures SoA = State-of-the-Art PACE = Physical+Adversarial AI+Cyber+EW

DoD cybersecurity operational test and evaluation ^[1]

[1] DoD Cybersecurity Test and Evaluation Guidebook, Version 2.0, Change 1. February 2020.

DoD cybersecurity operational test and evaluation ^[1]

DoD AI security operational test and evaluation

An AI red team employing SoA PACE techniques and red-teaming tools to execute novel AI kill chain TTPs

DoD = Department of Defense AI = Artificial Intelligence SoA = State-of-the-Art PACE = Physical+Adversarial AI+Cyber+Electronic Warfare TTP = Tactics, Techniques and Procedures

[1] DoD Cybersecurity Test and Evaluation Guidebook, Version 2.0, Change 1. February 2020.

Major gaps in the AI red-teaming ecosystem that DARPA can enable

AI kill chains needed for use by future DoD AI red teams

DoD = Department of Defense AI = Artificial Intelligence TTPs = Tactics, Techniques, and Procedures

AI red teams must consider the novel class of intersectional AI and cyber/EW attack vectors for AI attack effects

Hypothesis: An AI red team employing SoA PACE techniques and red teaming tools to execute novel AI kill chain TTPs can operationally assess AI-enabled battlefield systems more effectively than a team without them

An agent prescribes which techniques and tools to employ rather than randomly selecting vectors/effects

AI = Artificial Intelligence AAI = Adversarial AI EW = Electronic Warfare ASUT = AI-enabled System-under-Test TTP = Tactics, Techniques and Procedures

- SABER-OpX: A series of high-fidelity operational exercises to iteratively assess (red-team) already developed AI-enabled systems in the settings where they will be deployed
- Develop a "model" for an operational AI red-teaming process
- Evaluate the discriminative AI paradigm to scale across "battlefield" systems
- Emphasize AI/cyber/EW attack vectors that delivers AI attack effects
- Develop operationally-guided TTPs to create and execute AI kill chains
- Continuously integrate emerging techniques/tools into the AI red team
- Estimate operational security risk of deployed AI-enabled systems
- Determine and investigate the trade-space of AI red team effectiveness

Notional physically realizable adversarial attacks on infrared AI-enabled aided target recognition system

Blue team goal:

Enable autonomous UGV to navigate terrain to a goal

- Experiments evaluate AI-enabled autonomy in increasingly difficult operating environments
- Main metric: Increase average autonomous speed

AI red team goal: Disable autonomous UGV to navigate terrain to a goal

- AI red team will manipulate the autonomy by using PACE techniques/tools to degrade the AI-based perception
- Main metric: Reduce average autonomous speed

SABER-OpX #2: AI-enabled Autonomous Drone (Rotary-Wing)

Blue team goal:

Develop drone autonomy to harden system against EW

- Ongoing work to rapidly harden drone platforms to counter adversary EW effects, integrating AI to enable a drone to carry out its mission without direct human control
- Main metric: Success rate of autonomy to complete objective (recon mission, strike mission)

AI red team goal: Degrade autonomy to reduce drone system effectiveness

(left): https://www.skydio.com/x10 (right top): https://www.skydio.com/skydio-2-plus-enterprise (right bottom): https://medium.com/skydio/skydio-autonomy-a-new-age-of-drone-intelligence-12346111b2f1

- AI red team will manipulate the autonomy by using PACE techniques/tools to degrade the AI-based perception
- Main metric: Failure rate of autonomy to complete objective (recon mission, strike mission)

The ASUTs represent variety among the operational characteristics that impact counter-AI effects in near-term discriminative AI, implying that AI kill chain TTPs over the ASUTs should be applicable to new battlefield systems

Operational Characteristics	SABER-OpX #1	SABER-OpX #2	Joint Coverage?		
View Perspective	Ground	Air	\checkmark		
Illumination	Ambient+Self	Ambient Only	\checkmark		
Sensor Modalities	Optical, Thermal, GPS, Lidar, RADAR, NDVI	Optical, optional thermal	\checkmark		
Motion Dynamics	2D, constrained	3D, unconstrained	\checkmark		
Effect Delivery	Predictable Paths	Omnidirectional approach window	✓		

SABER program metric

Core program metric: AI red team effectiveness

 $R = 1/(\theta_P P + \theta_T T + \theta_C C)$

- *P* Performance degradation of Blue system given AI attack effect
 - $P = \frac{N}{P}$, see chart per SABER-OpX
- T Time (normalized) to generate the AI attack effect
- *C* Cost (normalized) to generate the AI attack effect
- θ_i Importance given to each criteria (sums to one)

^{1/}R is a proxy measure for operational security risk

Each SABER-OpX will be a 9-month sprint, with experiments occurring at the end of months 1, 3, 6, and 9

For each SABER-OpX, we aim to improve *R*:

- Experiment 1: establish baseline, R_{base}
- Experiment 2: 10% increase above *R*_{base} (goal)
- Experiment 3: 20% increase above *R*_{base} (metric)
- Experiment 4: 40% increase above *R*_{base} (goal)

SABER-OpX experiment metrics

SABER-OpX	System Metric (B)	Negative System Metric (N)					
OpX #1	Maximize average autonomy speed	Minimize average autonomy speed					
OpX #2	Mission success rate	Mission failure rate					

B – Measure of performance of the Blue system *without* red degradation N – Measure of performance of the Blue system *with* red degradation

SABER program schedule

	Q1		Q2			Q3			Q4			Q5			Q6			Q7			Q8			
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
		SABEI prepar	R-OpX rations	;	SABER-OpX #1 AI red-t							aming SAB			SAE	BER-OpX #2 AI red-teaming							Wrap-Up	
111	Survey, assess, select and integrate techniques and tools															Ref	fine olkit							
	Employ selected and integrated techniques and tools																							
	SABER-OpX #1 preparations and experimentation support																							
112	SABER-OpX #2 preparations and experimentation support																							
TT3		SABEF prepar	R-OpX rations	;		SAB	ER-OpX #1 AI red-teaming							SABER-OpX #2 AI red-tea					-tean	ning	I g Wrap-Up			
					Cour	nter-A	I R&D	, deve	elop no	ovel A	I kill c	hain T	TPs a	nd em	iploy s	electe	d/inte	grated	l techr	niques	and t	ools		
	Compile results, document TTPs, and estimate operational AI security ris													ks										

- ▲ = SABER OpX Experiment
- ▲ = SABER OpX Metric Evaluation

AI = Artificial Intelligence TTPs = Tactical, Techniques, Procedures

www.darpa.mil