

An Artificial Coevolutionary Framework for Adversarial AI

Jamal Toutouh
toutouh@mit.edu

ALFA

<http://groups.csail.mit.edu/ALFA>
unamay@csail.mit.edu



ALFA Group 2018-2019

Core Members



Una-May O'Reilly



Erik Hemberg



Abdullah Al-Dujaili



Jamal Toutouh

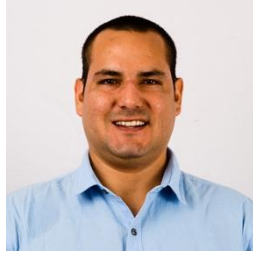


Nicole Hoffman

PhD



Shashank Srikant



Miguel Paredes

Research Assoc.



Edilberto Amorim



Jonathan Rubin

UROP



Milka Piszczek



Michal Shlapentokh-Rothman

MEng



Li Wang



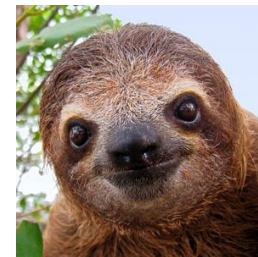
Jonathan Kelly



Ayesha Bajwa



Andrew Zhang



Linda Zhang

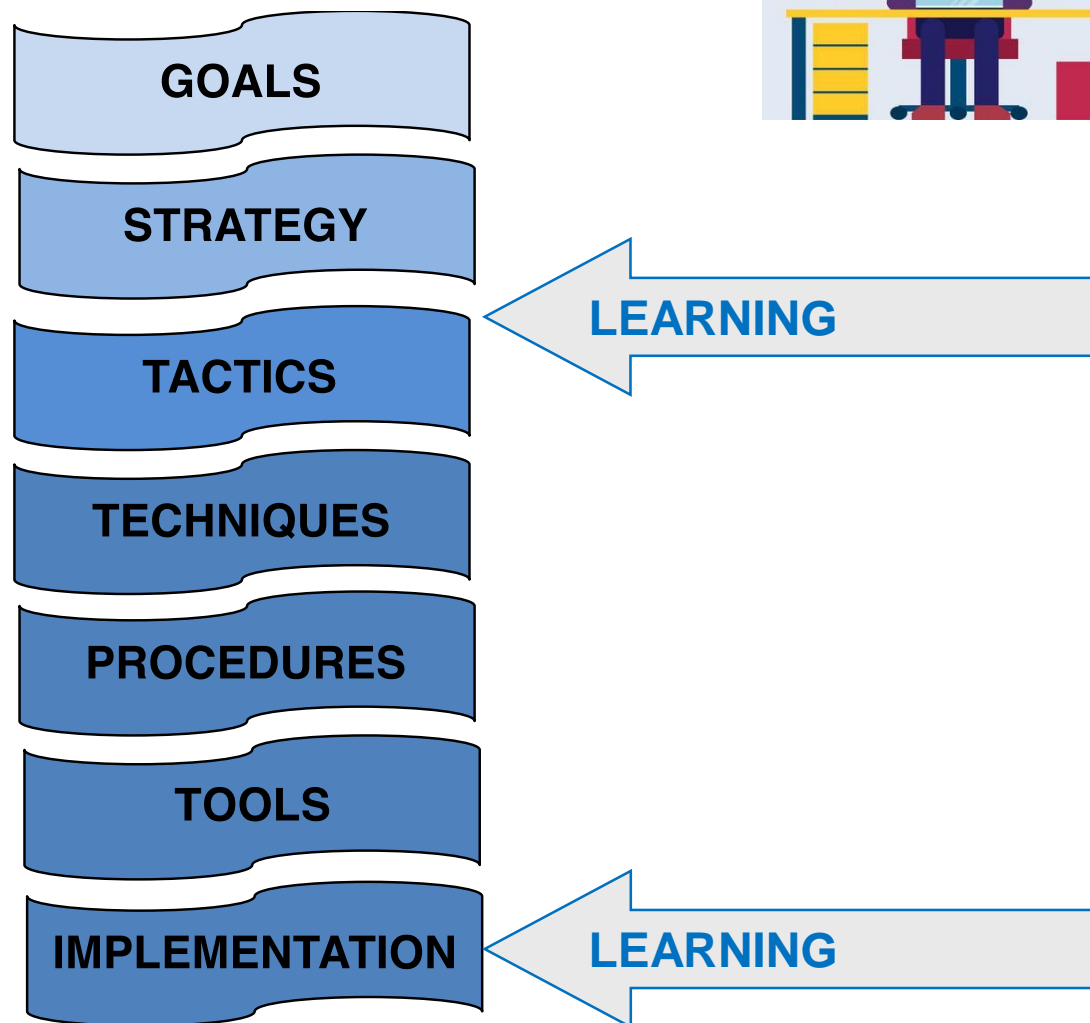
Grad



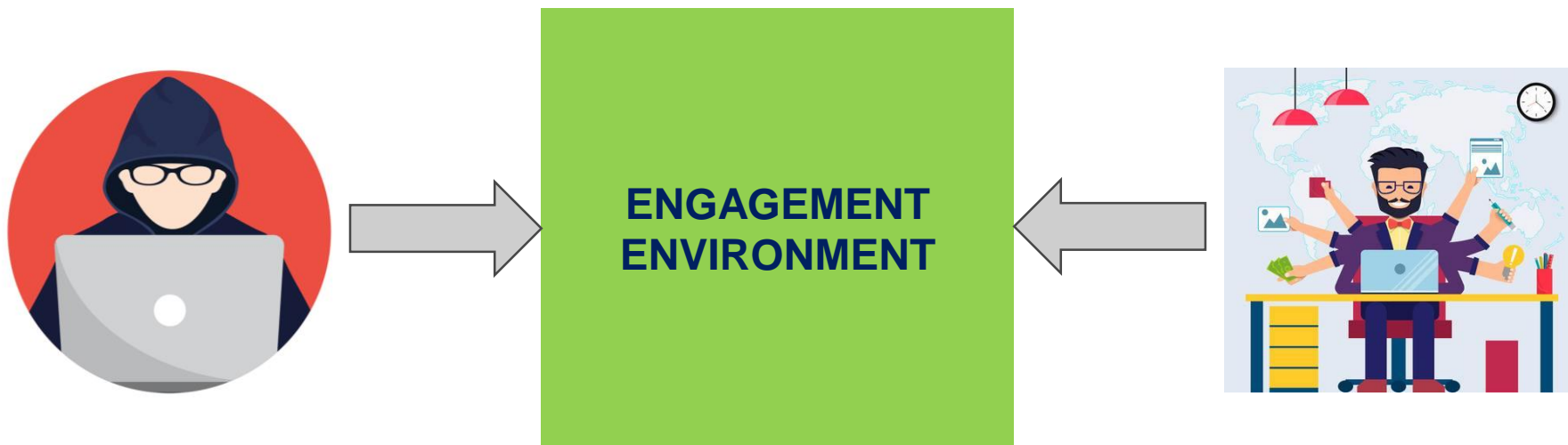
Saeyoung Rho (TPP)

Agenda

- **Adversarial Engagements and Arms Races**
- **Network Security Arms Races**
 - **RIVALS framework**
 - » **RIVALS: Robustness vs Denial**
 - » **AVAIL: Isolation vs Contagion**
 - » **DARK Horse and ADHD: Deception vs reconnaissance**
 - » **Acknowledgments:**
 - **Funding**
 - ❖ Member companies supporting Cybersecurity@CSAIL
 - ❖ DARPA XD3
 - ❖ MIT Lincoln Labs
 - **Work**
 - ❖ Members of the ALFA group and collaborators



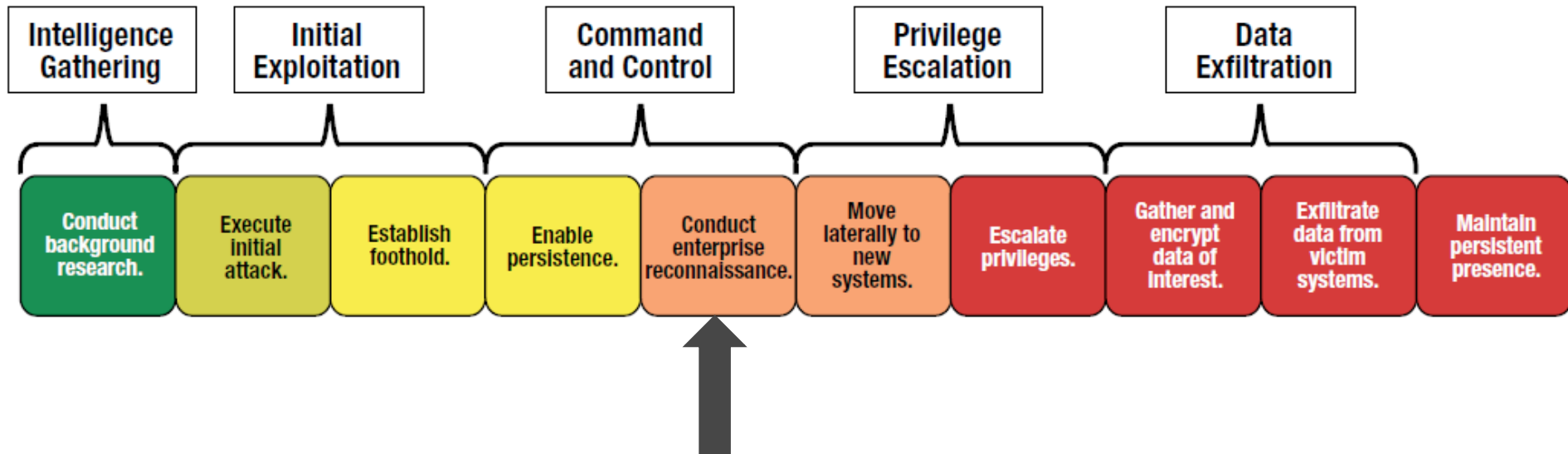
RIVALS



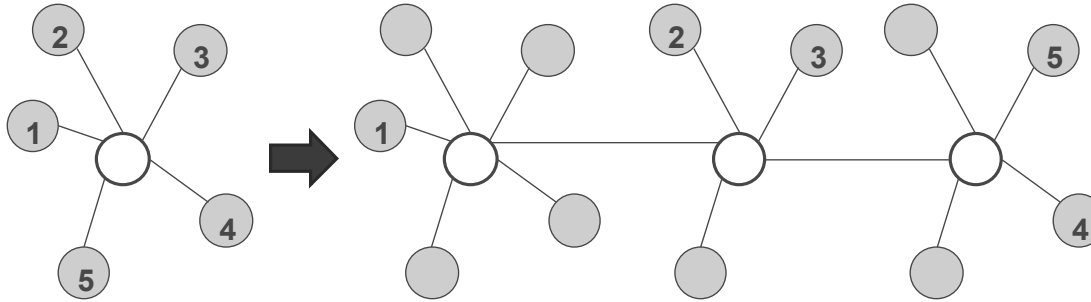
RIVALS helps the defense **anticipate** the attack strategies given a defensive configuration (and mission)

RIVALS helps the defense consider arms races and
Design effective courses of action for the network to be resilient

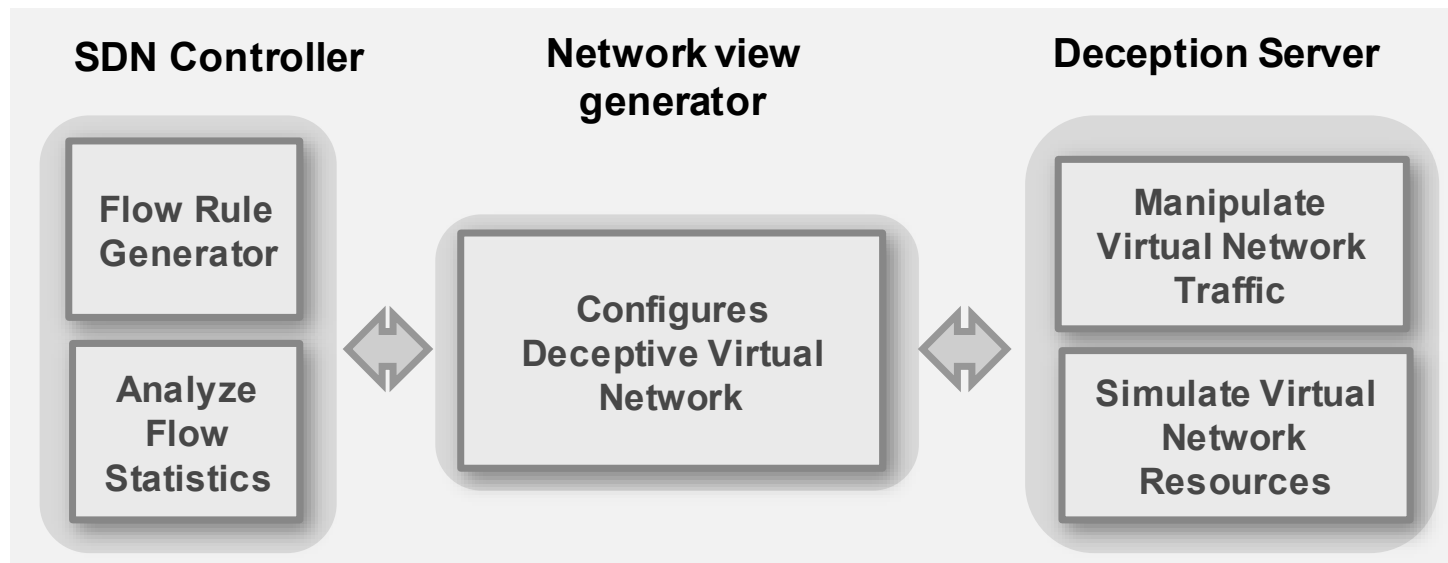
Advanced Persistent Threat Kill Chain



Deceptive Defense With Honeypots



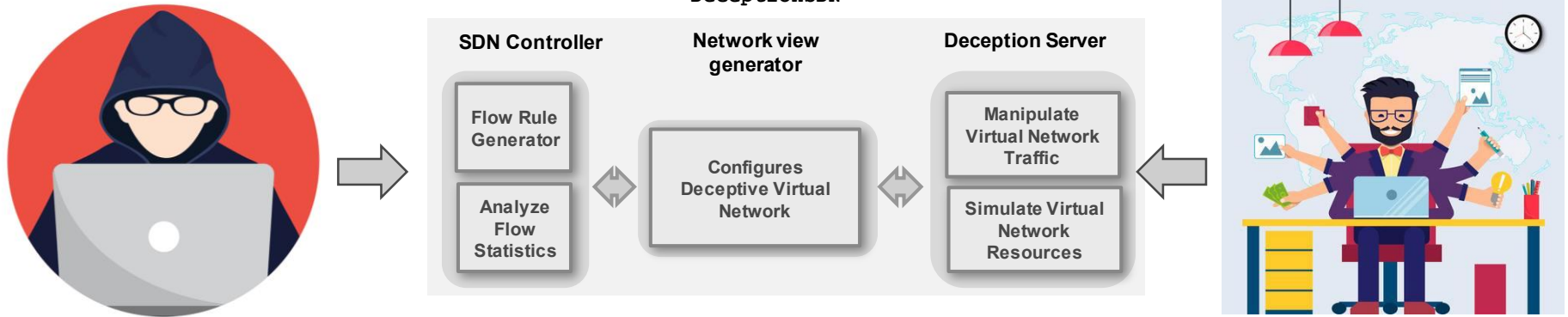
DeceptionSDN



Achleitner et al.

Engagement

CONFLICTING OBJECTIVES



MEASUREMENTS

d : time for defense to detect a scan (sec.)

t : time to run the scan (sec.)

n : number of scan detections by defender

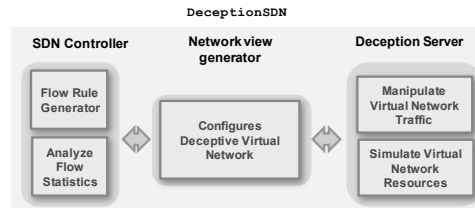
h/H : ratio of real nodes that were discovered to total real nodes.

Evaluate using Mininet

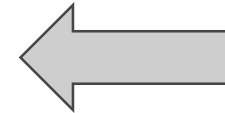
Adversarial Behaviors



APT
SCANNING
TACTICS

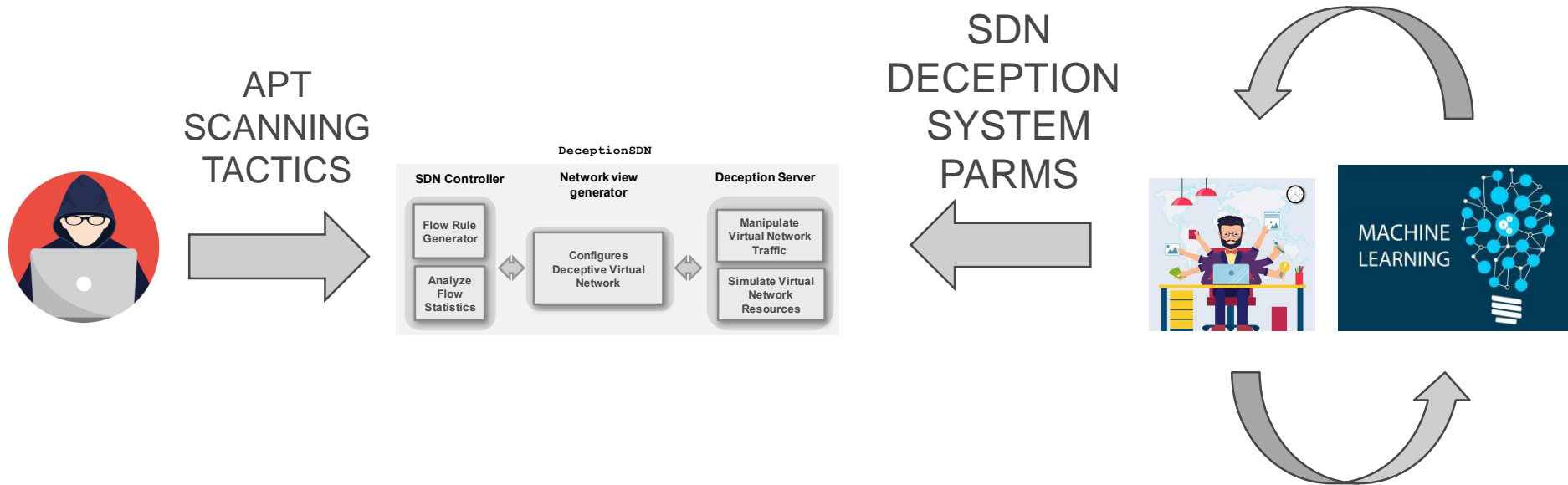


SDN
DECEPTION
SYSTEM
PARMS



Actor	Configuration	Range
Attacker	NMAP IP scan batch size	[10, 50, 100, 200]
	Total number of IPs to scan	[200, 300, 400]
	IP address visit order	<i>random, local, sequential, local-seq.</i>
Defender	Number of real nodes (H)	[20, ..., 40]
	Min honeypots per subnet	[1, ..., 10]
	Max honeypots per subnet	[11, ..., 20]
	Number of subnets	[3, ..., 6]
	Real node distribution in subnets	<i>random, even, crowding</i>

Defensive Learning



Static Attack – Optimized Defense

Hypothesis: Good defense has more honeypots, subnets and real hosts with even distribution

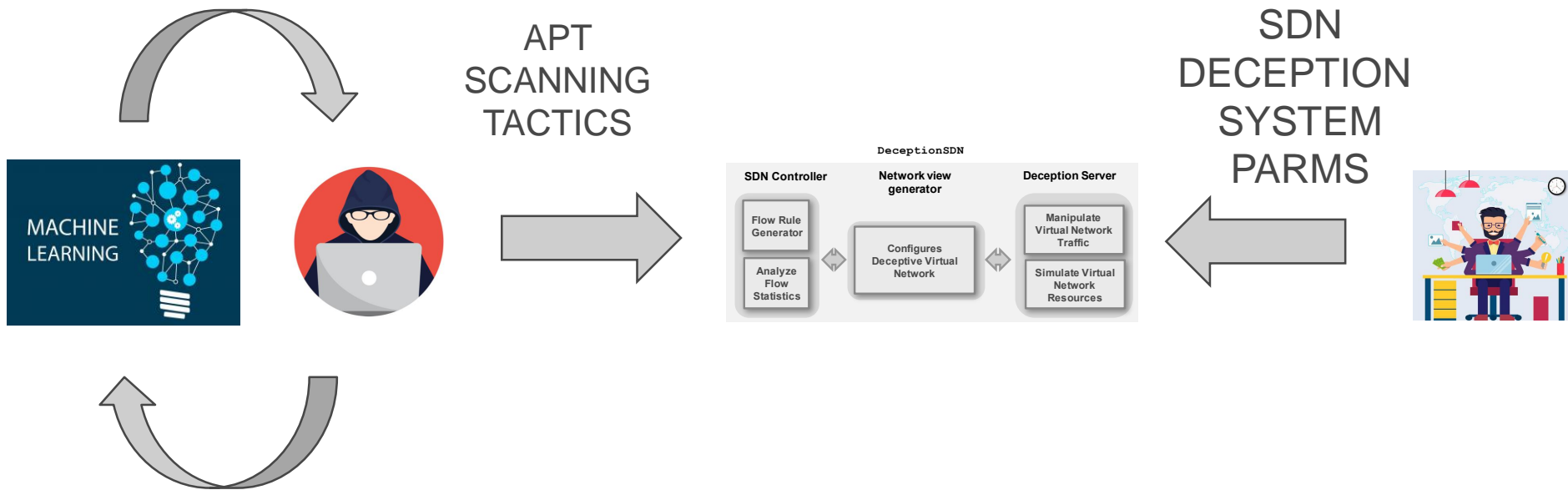
Results:

- More difficult to detect smaller NMAP batch sizes
 - Fitness function rewards discovering more real host less than the penalty of being detected: smaller scans do better
- Defense against an attacker that scans with local preference is the most difficult
- Expected real behavior of attackers is to start scan their local subnet

Possible recommendation: create subnets for DHCP leases where real hosts are in a different subnet

Visit Order	Batch Size	Num. IPs	N Real Nodes	Real Node Dist	Subnets	Min-Max HP	Nodes Disc.	Detected Scans	1 st Detection(s)	HPs
random	100	500	17	crowding	6	6-19	17	19	18.89	98
local	100	500	17	crowding	4	10-18	17	25	20.47	66
sequential	100	500	20	random	6	6-20	9	26	18.72	77
local-seq.	100	500	19	crowding	6	9-12	19	24	16.66	74
random	100	400	17	crowding	6	8-18	17	21	18.50	104
local	100	400	20	crowding	6	6-16	20	18	30.84	77
sequential	100	400	14	even	6	8-20	3	24	16.59	107
local-seq.	100	400	20	crowding	6	10-15	17	18	27.11	106
random	5	400	19	crowding	6	7-18	9	9	34.62	96
local	5	400	18	random	5	8-13	4	10	32.60	77
sequential	5	400	11	crowding	5	3-18	3	8	42.55	66
local-seq.	5	400	15	even	6	4-17	3	4	59.04	78
random	5	200	13	random	6	4-16	6	7	28.61	60
local	5	200	19	crowding	5	5-18	11	8	30.53	85
sequential	5	200	11	random	5	8-20	2	7	46.64	78
local-seq.	5	200	17	crowding	6	1-20	9	6	18.40	42
random	10	200	15	random	5	7-19	7	5	34.63	61
local	10	200	13	even	4	6-14	4	8	34.57	37
sequential	10	200	17	even	6	8-20	4	12	14.37	87
local-seq.	10	200	15	crowding	5	7-13	11	8	26.45	60

Attack Learning



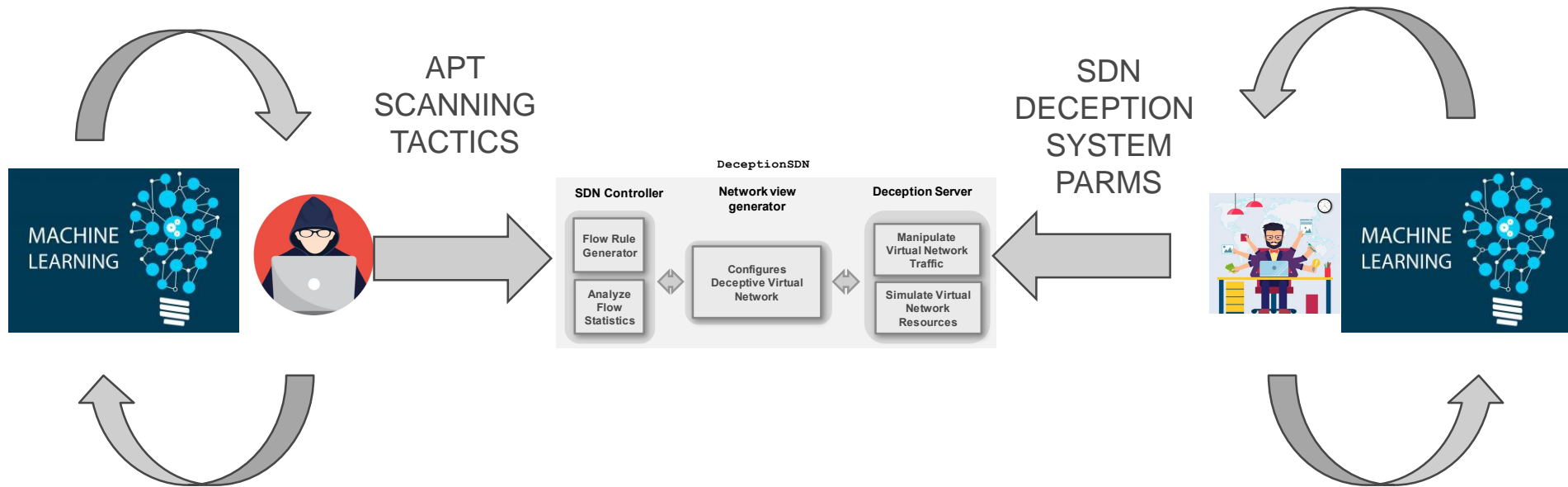
Static Defense – Optimized Attack

Results:

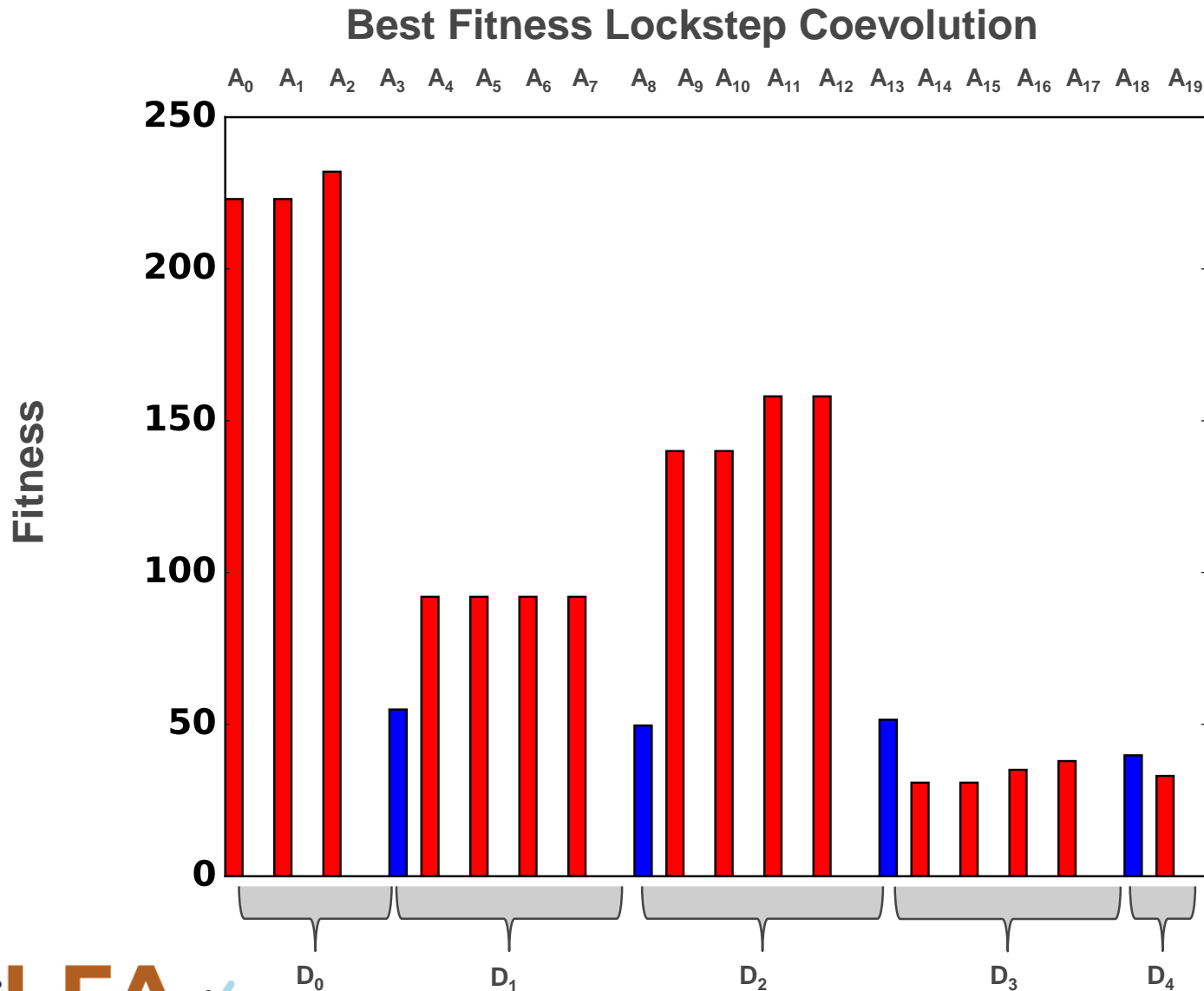
- Difficult to attack many honeypots and subnets.
 - Easier with crowded distribution of real hosts, large reward when that subnet is scanned (similar for defender when avoiding)
- Points to adopting high risk- and high reward tactic

N Real Nodes	Real Node Dist	Subnets	Min-Max HP	Visit Order	Batch Size	Num. IPs	Nodes Disc.	Detected Scans	1 st Detection(s)	HPs
20	even	4	1-10	local	25	300	8	11	22.62	26
20	crowding	4	1-10	sequential	5	400	12	2	42.74	14
20	random	4	1-10	local	10	300	4	3	26.61	19
20	even	4	10-20	local	25	200	5	19	22.54	86
20	crowding	4	10-20	sequential	50	300	20	19	16.47	73
20	random	4	10-20	sequential	10	400	8	12	22.43	76
20	even	10	1-10	sequential	5	400	3	2	24.52	54
20	crowding	10	1-10	sequential	50	400	20	11	20.50	68
20	random	10	1-10	local	5	200	2	1	85.35	54
20	even	10	10-20	local	5	200	1	9	40.72	185
20	crowding	10	10-20	sequential	50	300	20	22	23.06	214
20	random	10	10-20	local	5	200	5	9	40.64	170
20	even	10	40-80	sequential	50	200	2	4	22.13	983
20	crowding	10	40-80	sequential	50	300	11	18	18.64	683
20	random	10	40-80	sequential	5	200	2	12	30.84	998
50	even	10	1-10	local-seq.	25	200	9	13	30.58	36
50	crowding	10	1-10	sequential	5	300	15	2	68.95	52
50	random	10	1-10	sequential	5	400	3	2	97.45	48
50	even	10	10-20	sequential	25	200	7	16	20.42	205
50	crowding	10	10-20	local-seq.	5	400	20	7	30.76	231
50	random	10	10-20	sequential	50	200	2	14	22.63	183
20	even	20	10-20	sequential	25	200	2	17	16.67	379
20	crowding	20	10-20	local	5	300	12	4	24.81	347
20	random	20	10-20	local-seq.	25	300	1	18	16.58	389

Coevolutionary Arms Race

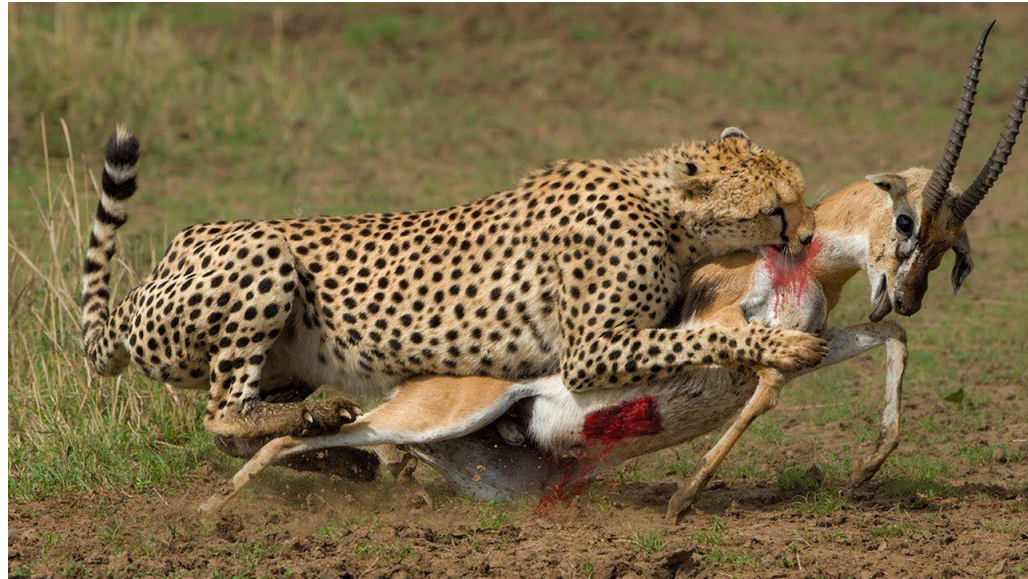


Coevolution of Scanning and Deception



Defender does significantly worse vs an evolving attacker, thus beware of static assumptions

Evolved for defense & attack

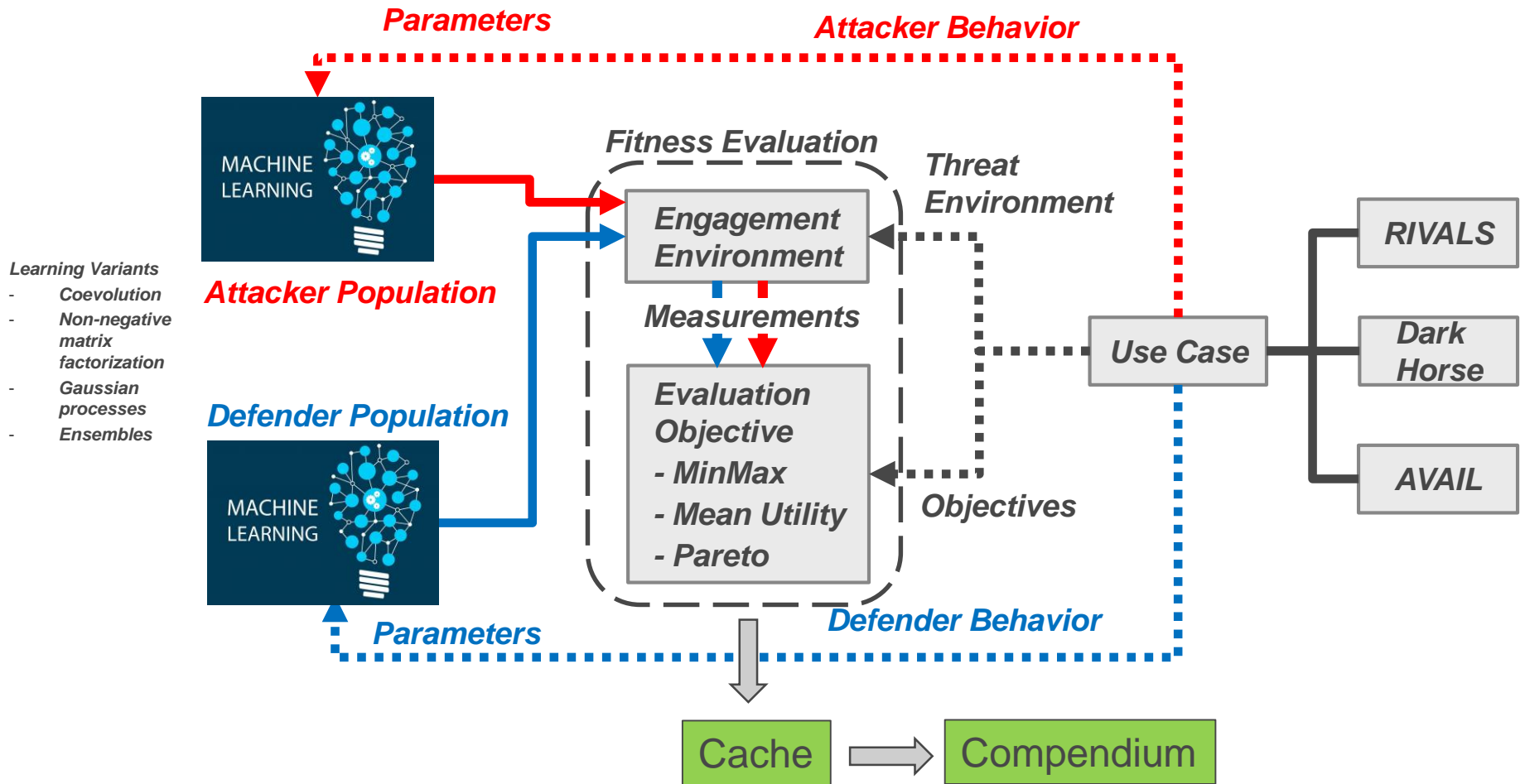


From Biological Coevolution Towards Adversarial AI Via Artificial Coevolution



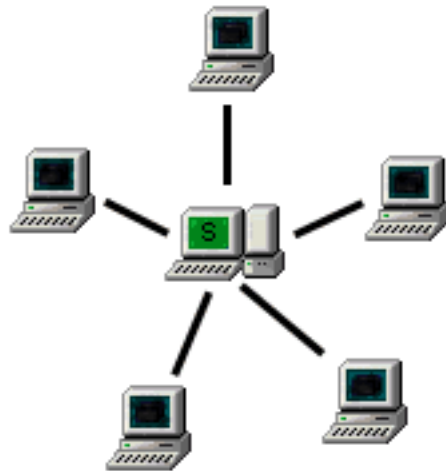
- **Biological arms races can provide adaptation**
 - **Can coevolution help to improve robustness in other adversarial settings?**
 - **Multiple comparisons can aid robustness and improve diversity**
 - **Help to anticipate**
 - **Replay the arms-race**

Adversarial AI Framework Concept

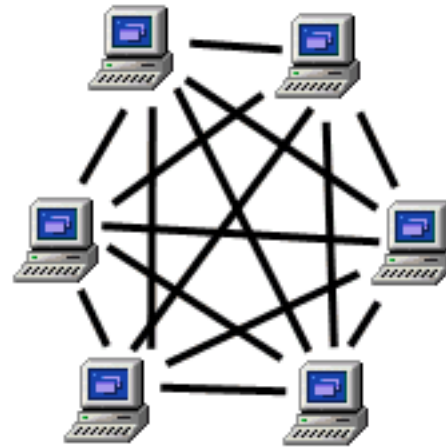


DDoS Network Defense

Server Based Network



Peer to Peer Network



https://proxy.duckduckgo.com/iu/?u=https%3A%2F%2Fcdn-images-1.medium.com%2Fmax%2F1600%2F1*8gG2DMHmnGJBlyDtR79rjA.png&f=1

RIVALS: Network Routing Problem

FITNESS

mission disruption
attacks in number and duration



$$f_a^L = \frac{1 - \text{mission_success}}{(n_attacks \cdot \text{total_duration}) + n_attacks}$$



FITNESS

mission completion
time and hops

$$f_d^L = \frac{\text{mission_success}}{\text{overall_time} \cdot n_hops}$$

Defender Objective: maximize

Attacker Objective: maximize

ATTACKER ACTIONS

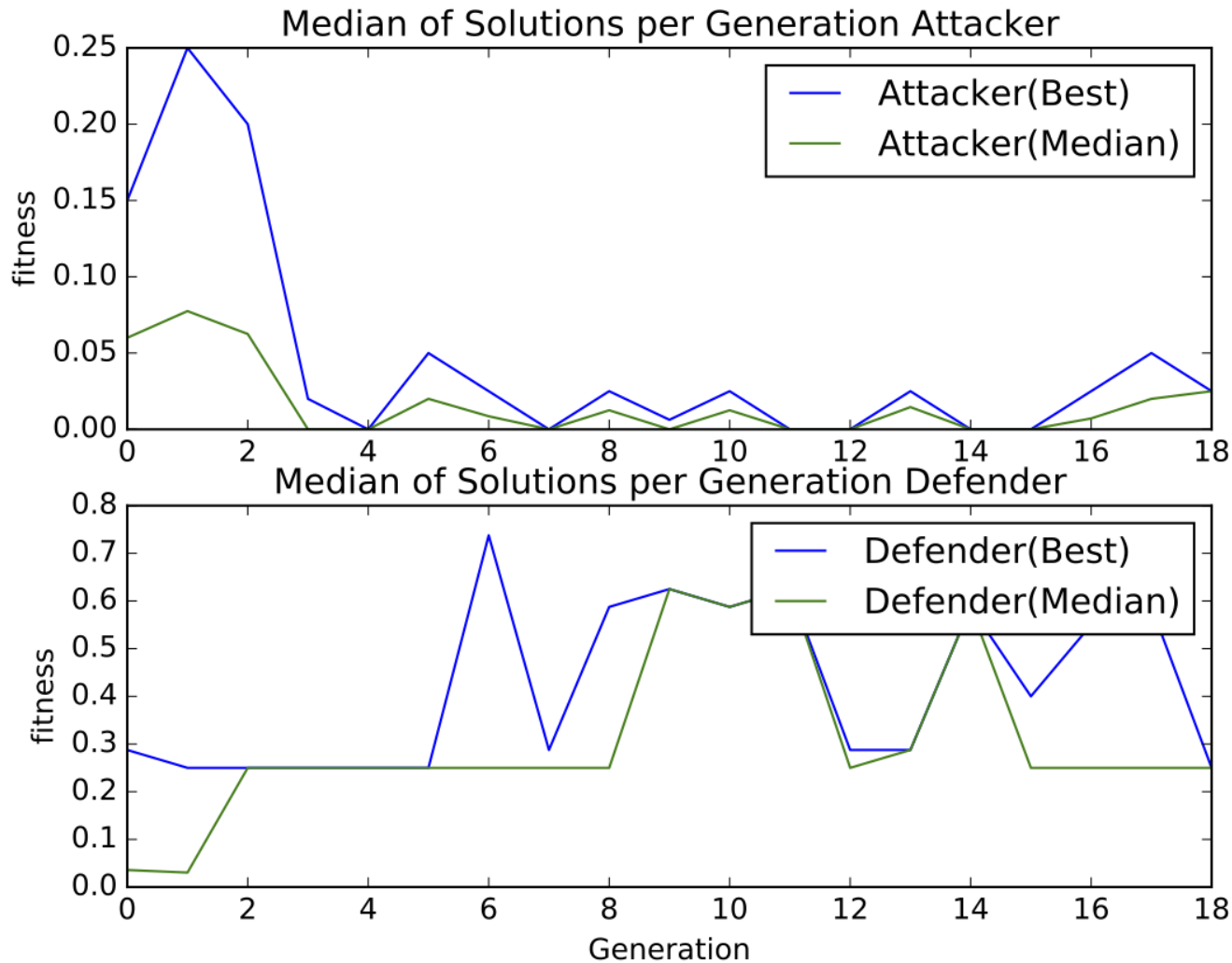
node, start time, end time
complete loss of node

DEFENDER ACTIONS

link flooding
shortest path
CHORD

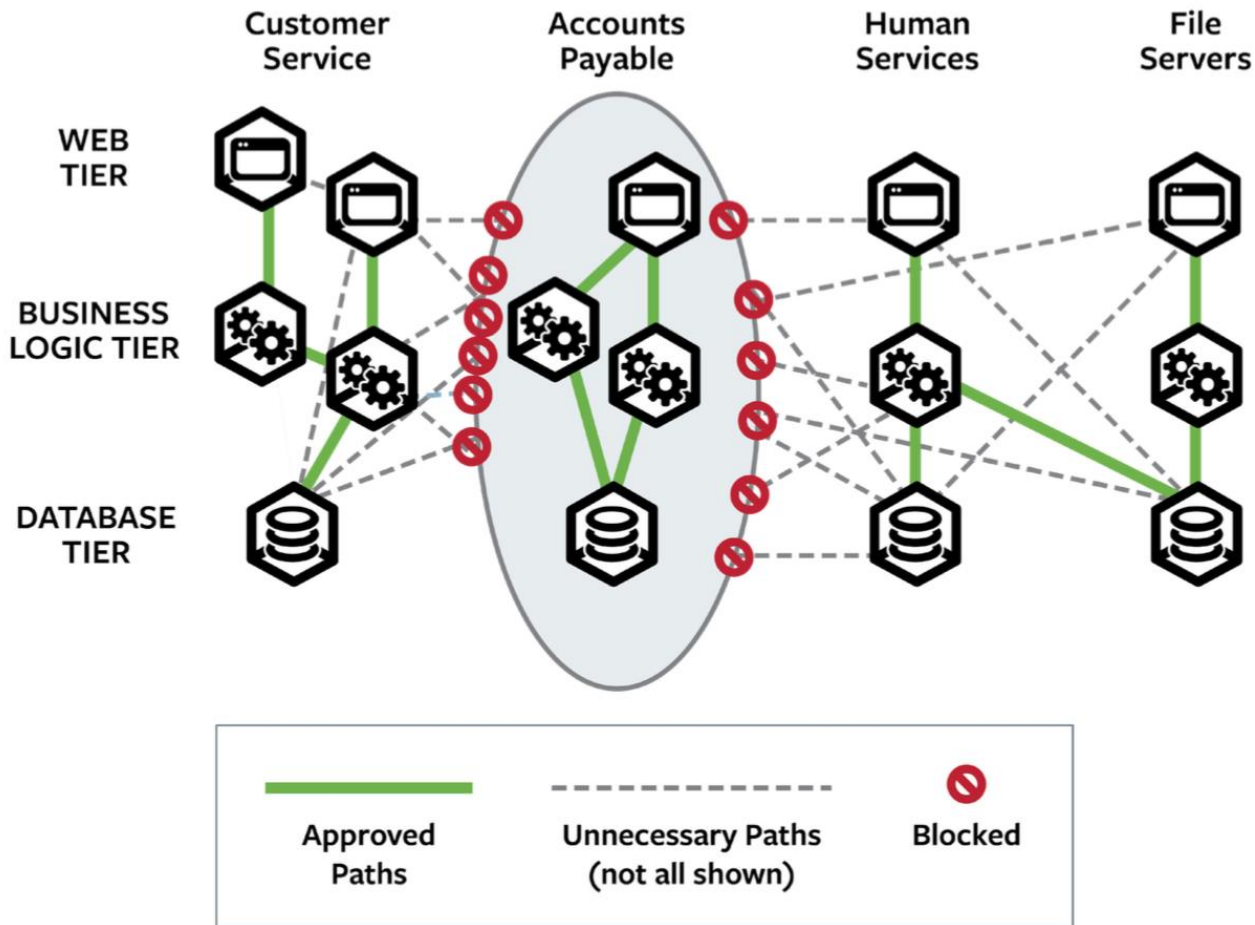


ALFA Sim of P2P



Defender finds an optimal solution

Network Segmentation



The Definitive Guide to Micro-Segmentation, John Friedman, CyberEdge Group

AVAIL: Enclaves vs Contagion

FITNESS

mission delay
budget remaining



FITNESS

mission delay
budget remaining

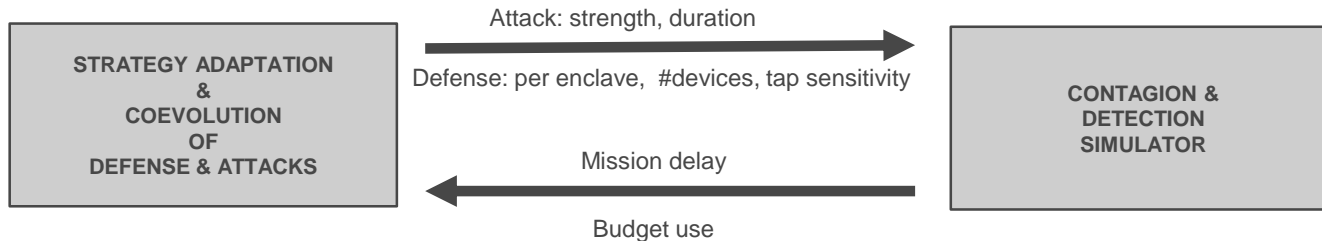
OBJECTIVE: minmax

ATTACKER ACTIONS

set strength and duration of attack
for each enclave

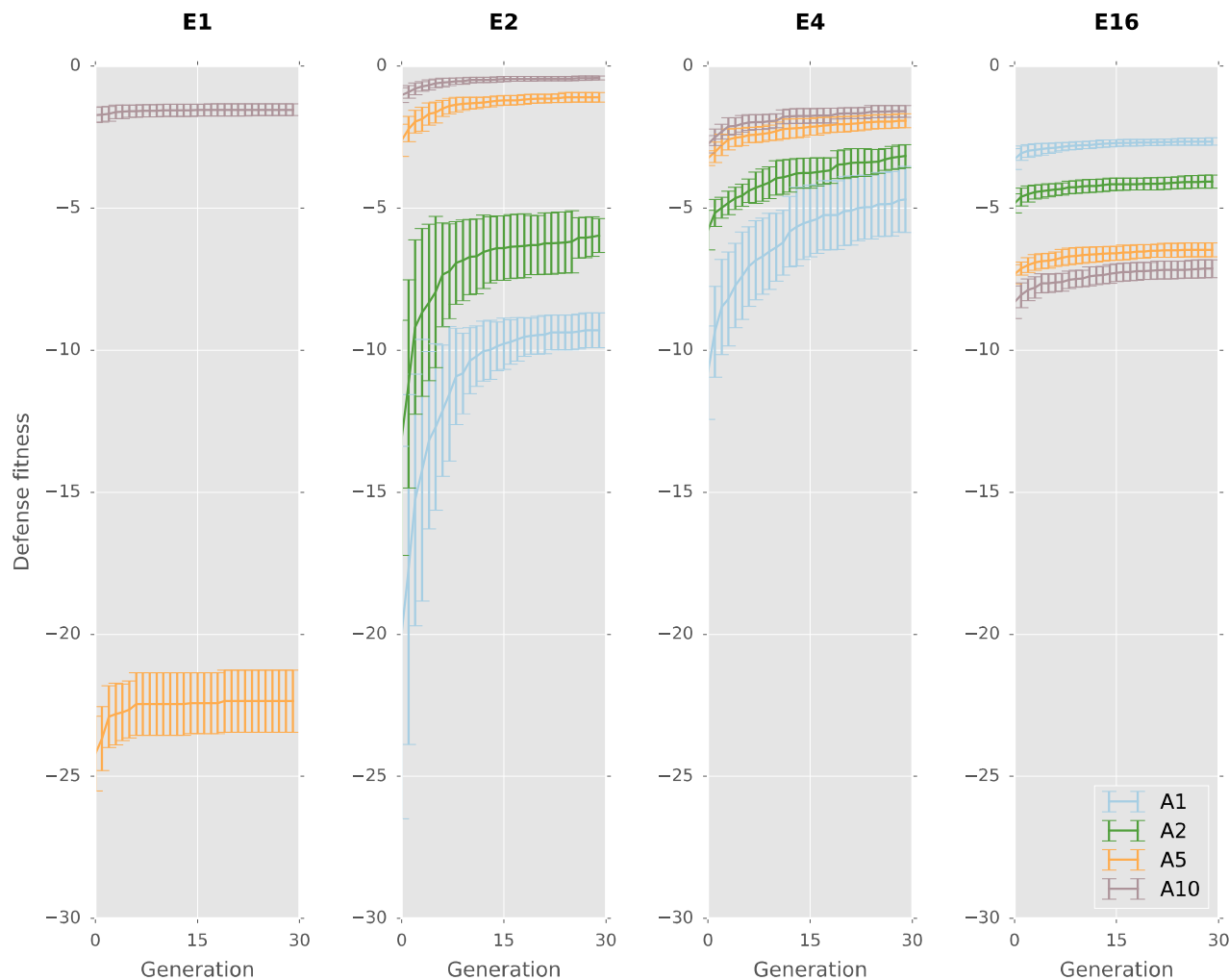
DEFENDER ACTIONS

set tap sensitivity and size
for each enclave



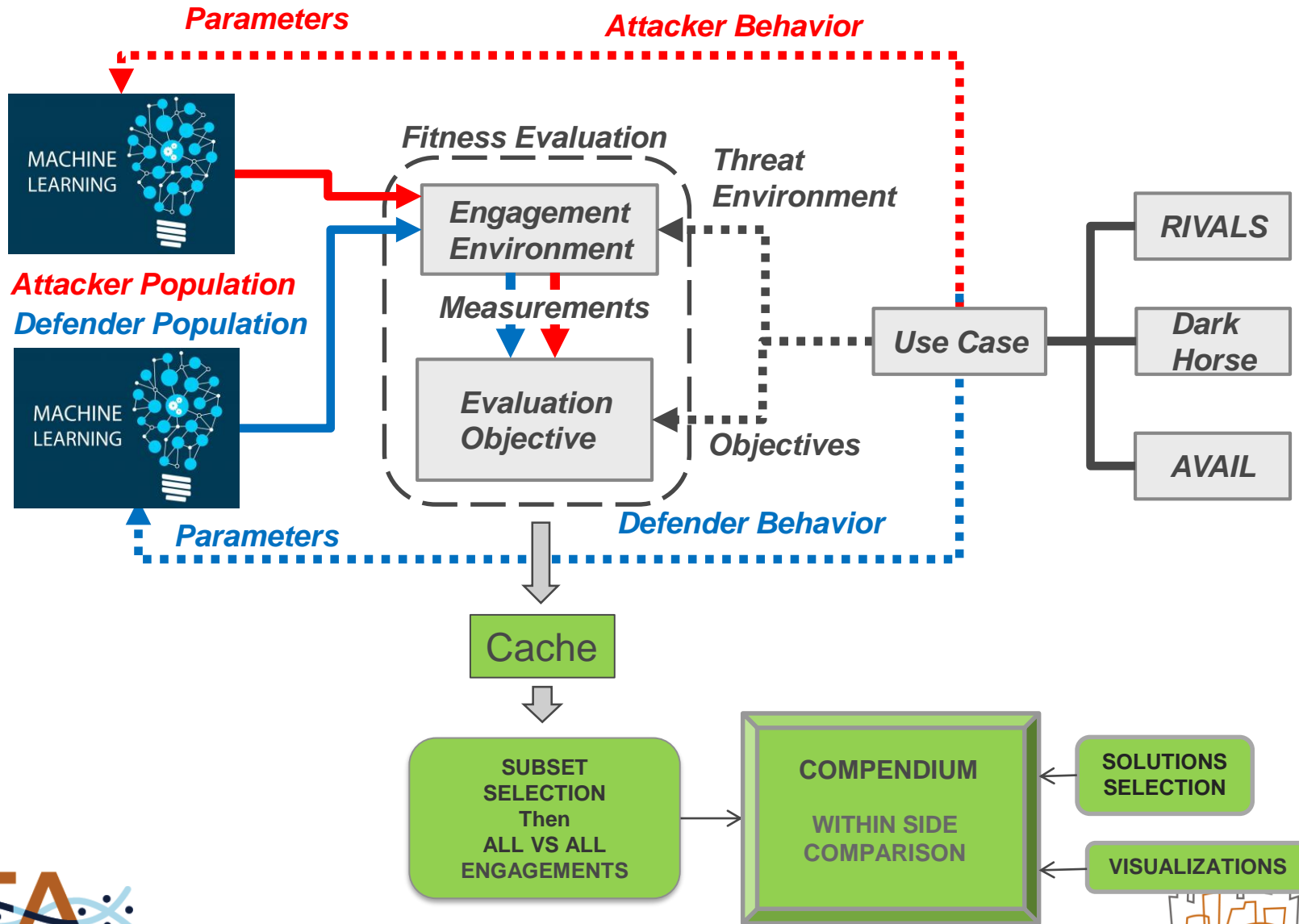
Realistic parameters from a
validated low level emulator

Evolve Defense



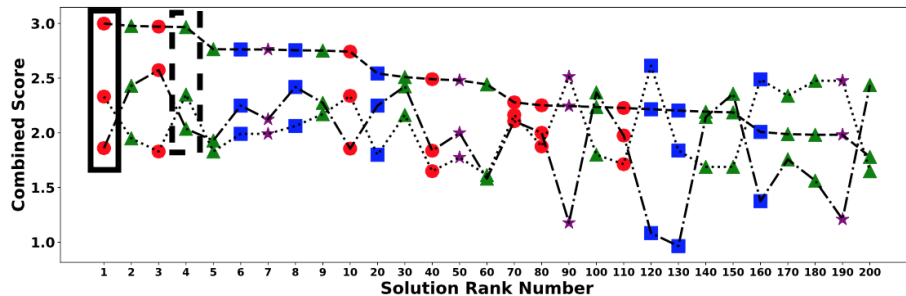
Defender learns sound network segmentation practices over generations

Compendium Analysis

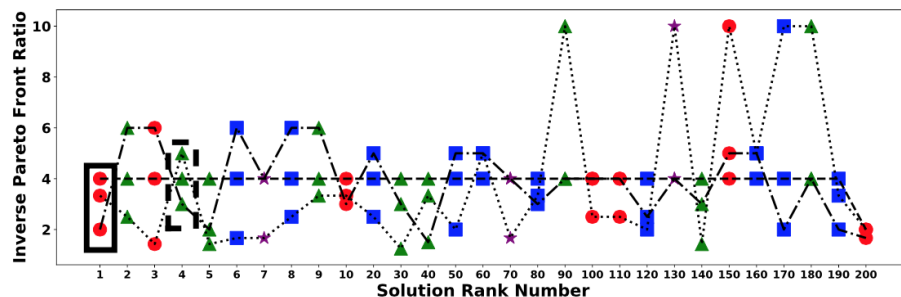


Attack Campaign Performance Comparison

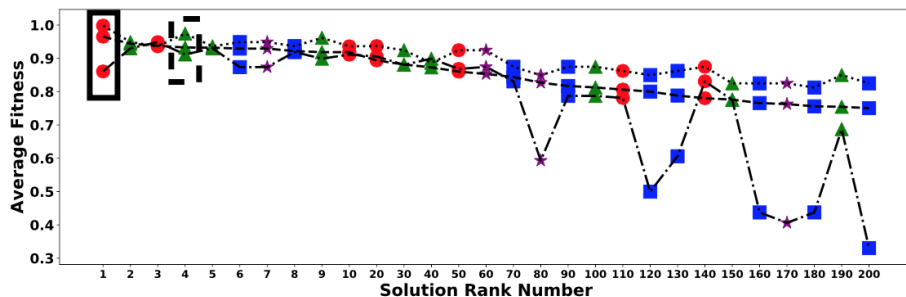
Different metrics and Ranking Schemes



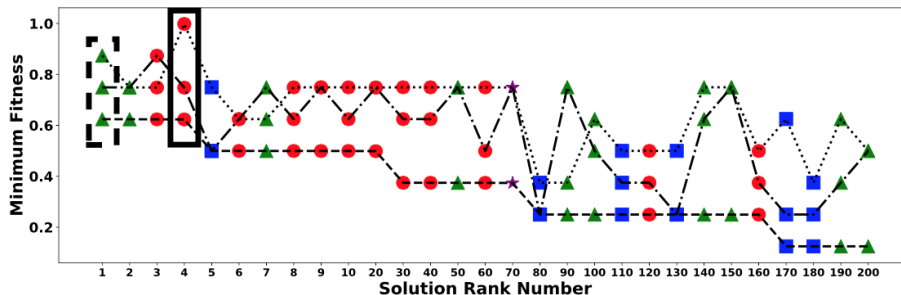
(a) Attacker CS ranking scheme.



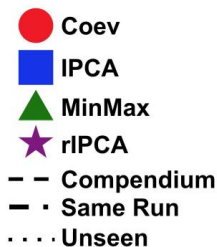
(b) Attacker PF ranking scheme.



(c) Attacker AF ranking scheme.

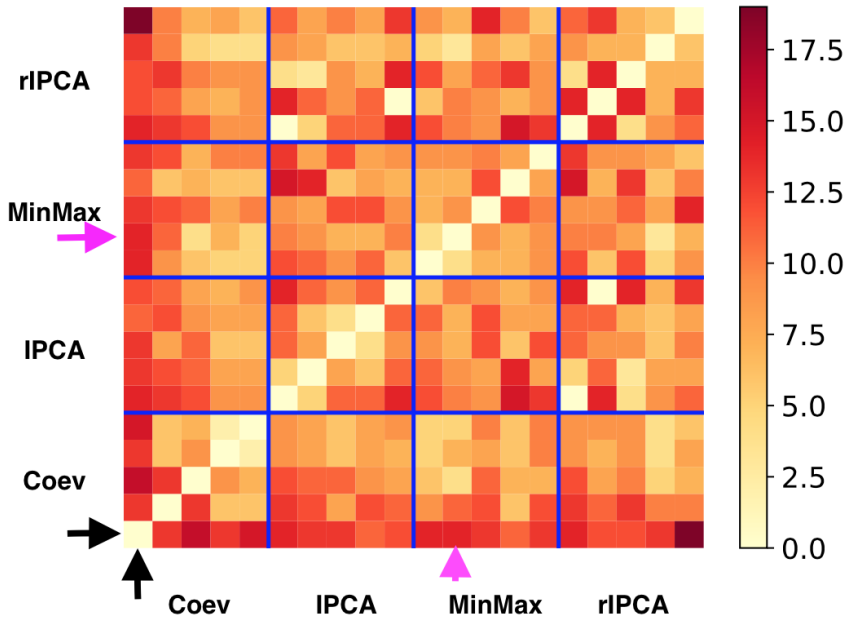


(d) Attacker MF ranking scheme.



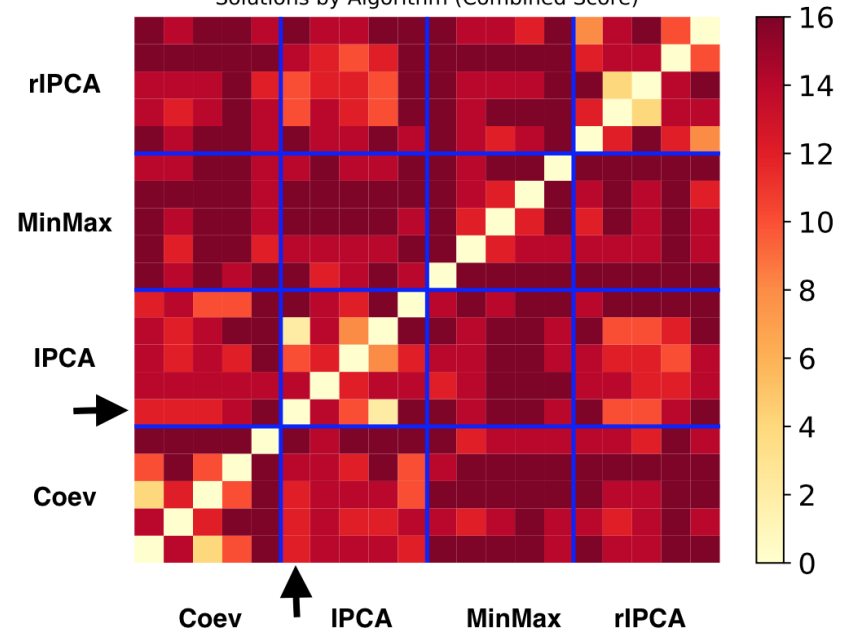
Attack Campaign Similarity

Phenotype Distances for Top Attacker Solutions by Algorithm (Combined Score)



(a) Attacker pairwise distance. Black arrow shows attacker selected by AF, PF, and CS ranking scheme. The pink arrow shows the attacker selected by the MF ranking scheme.

Phenotype Distances for Top Defender Solutions by Algorithm (Combined Score)



(b) Defender pairwise distance. Black arrow shows defender selected by all ranking schemes.

Summary & Future Work

- **Adversarial Engagements and Arms Races**
- **Network Security Arms Races**
 - **RIVALS Adversarial AI framework**
 - » **RIVALS: Robustness vs Denial**
 - » **AVAIL: Isolation vs Contagion**
 - » **DARK Horse and ADHD: Deception vs reconnaissance**
- **Future**
 - **Validate, refine, and extend use cases**