

# Learning to Sail Dynamic Networks: The MARLIN Reinforcement Learning Framework for Congestion Control in Tactical Environments

Raffaele Galliera<sup>§,†</sup>, Mattia Zaccarini<sup>‡</sup>, Alessandro Morelli<sup>§</sup>  
Roberto Fronteddu<sup>§</sup>, Filippo Poltronieri<sup>‡</sup>, Niranjan Suri<sup>\*,§,†</sup>, Mauro  
Tortonesi<sup>‡</sup>

<sup>†</sup>The University of West Florida (UWF)

<sup>‡</sup>The University of Ferrara (UNIFE)

<sup>§</sup>Institute for Human and Machine Cognition (IHMC)

<sup>\*</sup>US Army DEVCOM Army Research Laboratory (ARL)



# OVERVIEW

Background

The Congestion Control Environment  
Challenges

Method

Experimental Setting

Results

Backup

# CONGESTION CONTROL

- ▶ The channel can saturate causing delays and re-transmissions
- ▶ Different heuristics are used (eg, **TCP Cubic**).
- ▶ **Congestion Window (CWND)**: control over bytes allowed to be **in-flight**.
- ▶ Heavy congestion can take a link to an impracticable state.

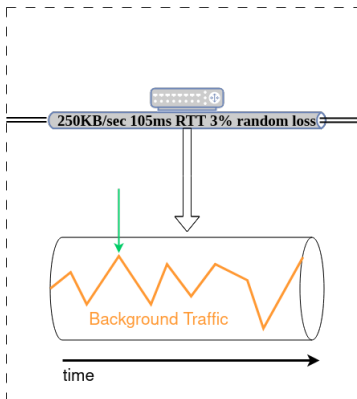
# STARTING POINT



Physical Training network.

R. Galliera, A. Morelli, R. Fronteddu and N. Suri, "MARLIN: Soft Actor-Critic based Reinforcement Learning for Congestion Control in Real Networks", NOMS 2023 IEEE/IFIP Network Operations and Management Symposium

# THE AGENT'S PERCEPTION



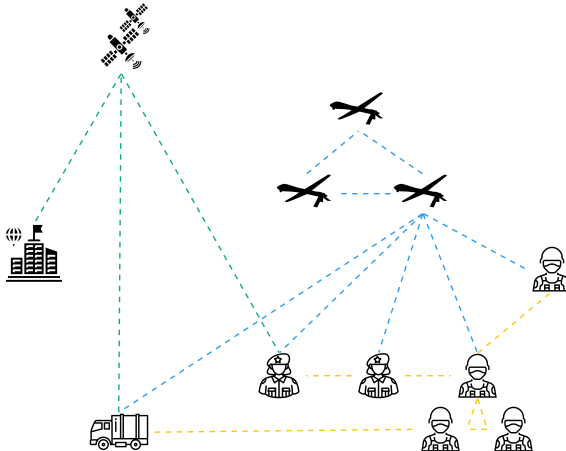
R. Galliera, A. Morelli, R. Fronteddu and N. Suri, "MARLIN: Soft Actor-Critic based Reinforcement Learning for Congestion Control in Real Networks", NOMS 2023 IEEE/IFIP Network Operations and Management Symposium



# MARLIN SUMMARY

- ▶ **Partnering protocol:** Mockets.
- ▶ **Non-blocking** communication.
- ▶ Third-party sources utilize the **shared link**
- ▶ **Learning algorithm:** Soft Actor-Critic (SAC).
- ▶ **Action space:**  $[-1, 1]$  tweaking the CWND.
- ▶ **Action-Observation History:** 10 steps.

# TACTICAL NETWORKING ENVIRONMENTS





# CHALLENGES

- ▶ **Designing networking scenarios** might become unfeasible.
- ▶ **Reproducibility.**
- ▶ **Reinforcement signals** in tactical networks.





# OBJECTIVE

Design a **flexible framework** to train agents in **dynamic** and **unreliable networking** scenarios.



# OVERVIEW

Background

Method

Framework Design

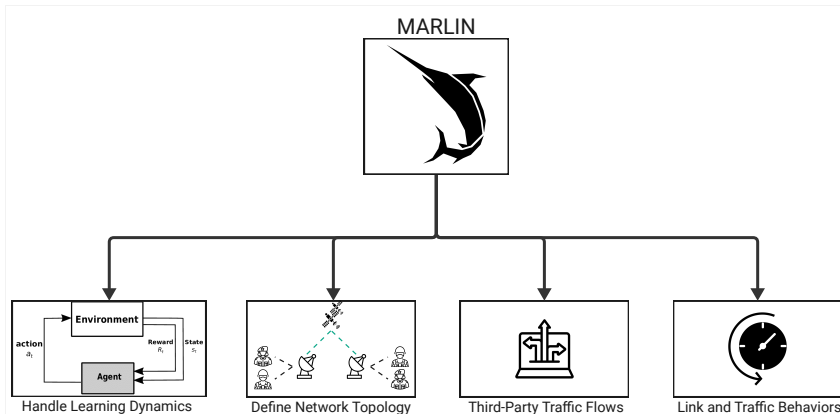
Rewarding in Tactical Networks

Experimental Setting

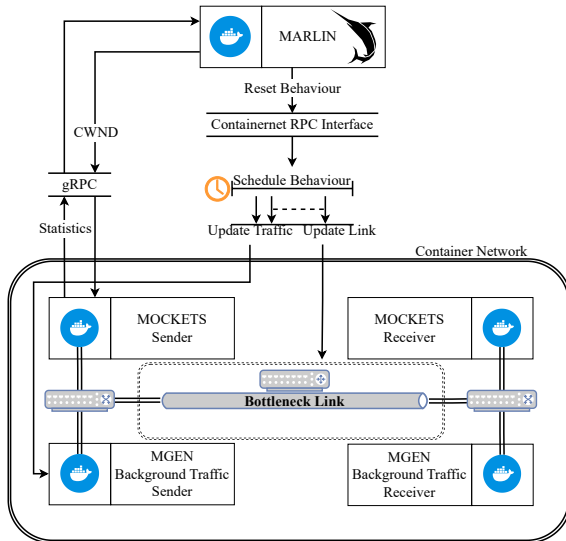
Results

Backup

# EXTENDING THE MARLIN FRAMEWORK



# SCHEDULING LINK BEHAVIORS



# REWARD

$$r_t = - \frac{\text{target}_t [1 + \text{retr} * (1 - \text{loss}_c)]}{\text{target}_t + \text{acked}_t^{\text{cumulative}}} \quad (1)$$



# OVERVIEW

Background

Method

Experimental Setting  
Networking Scenario  
Evaluating the Agent

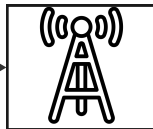
Results

Backup

# LINK TRANSITIONS



**SATCOM link**  
Bandwidth: 1Mb/s  
Delay: 500ms

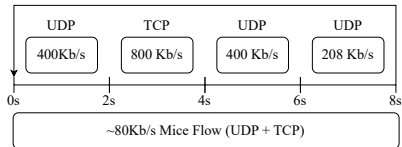


**UHF Radio link**  
Bandwidth: 256Kb/s  
Delay: 125ms

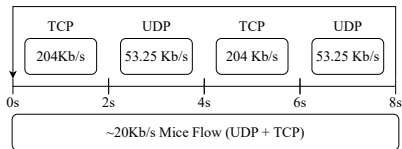
Major training details:

- ▶ 500K steps
- ▶ Fixed random packet loss (0% SATCOM, 3% UHF Radio)
- ▶ Link switch after 10 Seconds from the episode start

# THIRD-PARTY TRAFFIC BEHAVIORS



(a) Traffic during **SATCOM** link



(b) Traffic during **UHF** link

Flows generated with Multi-Generator (MGEN) Network Test Tool - U.S. Naval Research Laboratory





# TESTING SCENARIO

- ▶ **Objective:** 600KB payload transfer
- ▶ Varying UHF radio link random packet loss (0-3%)
- ▶ 100 testing episodes for each packet loss value
- ▶ Metrics used:
  - ▶ Transfer time (s)
  - ▶ Retransmissions
  - ▶ RTT Transition Impact

## RTT TRANSITION IMPACT (RTI)

$$RTI = \ln \left( \frac{\sum_{i=1}^m \frac{r_{tt_{i,max}}}{r_{tt_{i,nom}}}}{m} \right) \quad (2)$$

- ▶  $m$  link transitions
- ▶  $r_{tt_{i,max}}$  maximum  $r_{tt}$  detected during link  $i$
- ▶  $r_{tt_{i,nom}}$  nominal  $r_{tt}$  value during link  $i$



# OVERVIEW

Background

Method

Experimental Setting

Results

RTI

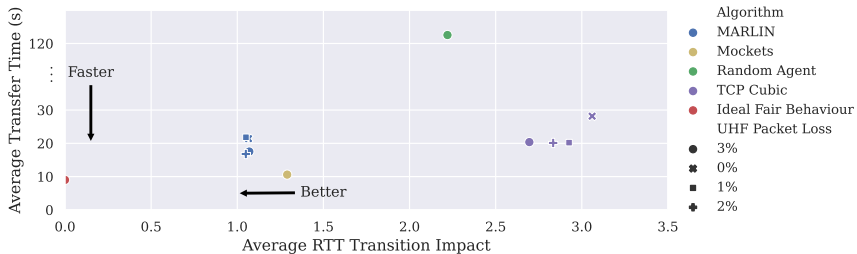
Retransmissions

Conclusion

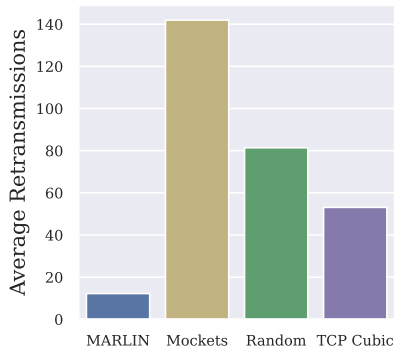
Backup



# AVERAGE TRANSFER TIME - RTI



# RETRANSMISSIONS

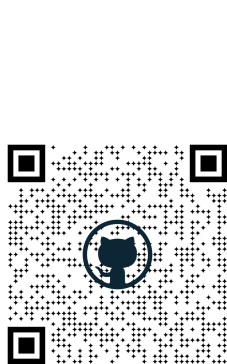


During experiments with 3% random packet loss set on the UHF link.



## CONCLUSION

- ▶ **A containerized RL environment** to train agents for CC.
- ▶ **Centralized control** of the entire training environment.
- ▶ **Programmable link behaviour.**
- ▶ **Retransmission-sensitive** rewards.
- ▶ **Competitive** trained policies.



GitHub Repository



# OVERVIEW

Background

Method

Experimental Setting

Results

Backup

# HYPERPARAMETERS

Hyperparameter	Value
Training steps	$5 \times 10^5$
History length	10
Training episode length	200
Learning rate	$3 \times 10^{-4}$
Buffer size	$2.5 \times 10^5$
Warm-up (learning starts)	$1 \times 10^4$ steps
Batch size	512
Tau	$5 \times 10^{-3}$
Gamma	0.99
Training Frequency	$1/\text{episode}$
Gradient Steps	-1 (same as episode length)
Entropy regularization coefficient	"auto" (Learned)
MLP policy hidden layers	[400, 300]

**Table:** Hyperparameters used in our experiment.





# STATE

	Feature	Description		Statistic
1	Current cwnd	Current cwnd	1	Last
2	KBs Sent	Amount of KB sent *	2	Mean
3	New KBs sent	Amount of KB acked *	3	STD
4	Acked KBs	Amount of KB acked *	4	Min
5	Packets sent	Packets sent *	5	Max
6	Retransmissions	Number of packets retransmitted *	6	EMA
7	Instantaneous Throughput	Throughput *	7	Difference from Previous
8	Instantaneous Goodput	Goodput *		
9	Unacked KBs	Amount of KBs in flight		
10	Last RTT	Last rtt detected *		
12	Min RTT	Min rtt *		
12	Max RTT	Max rtt *		
13	SRTT	Smoothed rtt *		
14	VAR RTT	rtt variance *		
		* During the last rtt timeframe		

Every feature has 7 nested statistics with a 10 observations history.



# IMPLEMENTATION STACK

- ▶ **Reinforcement Learning framework:** Stable Baselines 3.
- ▶ **Partnering Protocol:** Mockets.
- ▶ **Protocol-Agent Communication:** gRPC.
- ▶ **Network Emulation:** Containernet + RPYC.
- ▶ **Third-Party traffic:** MGEN