

RL-based Internet Congestion Control

Ashwin Shenai (180156), Astitva Chaudhary (180157)

EE698V: Advanced Topics in ML for Communication Networks

Fall 2022

Project Goals: Proposed

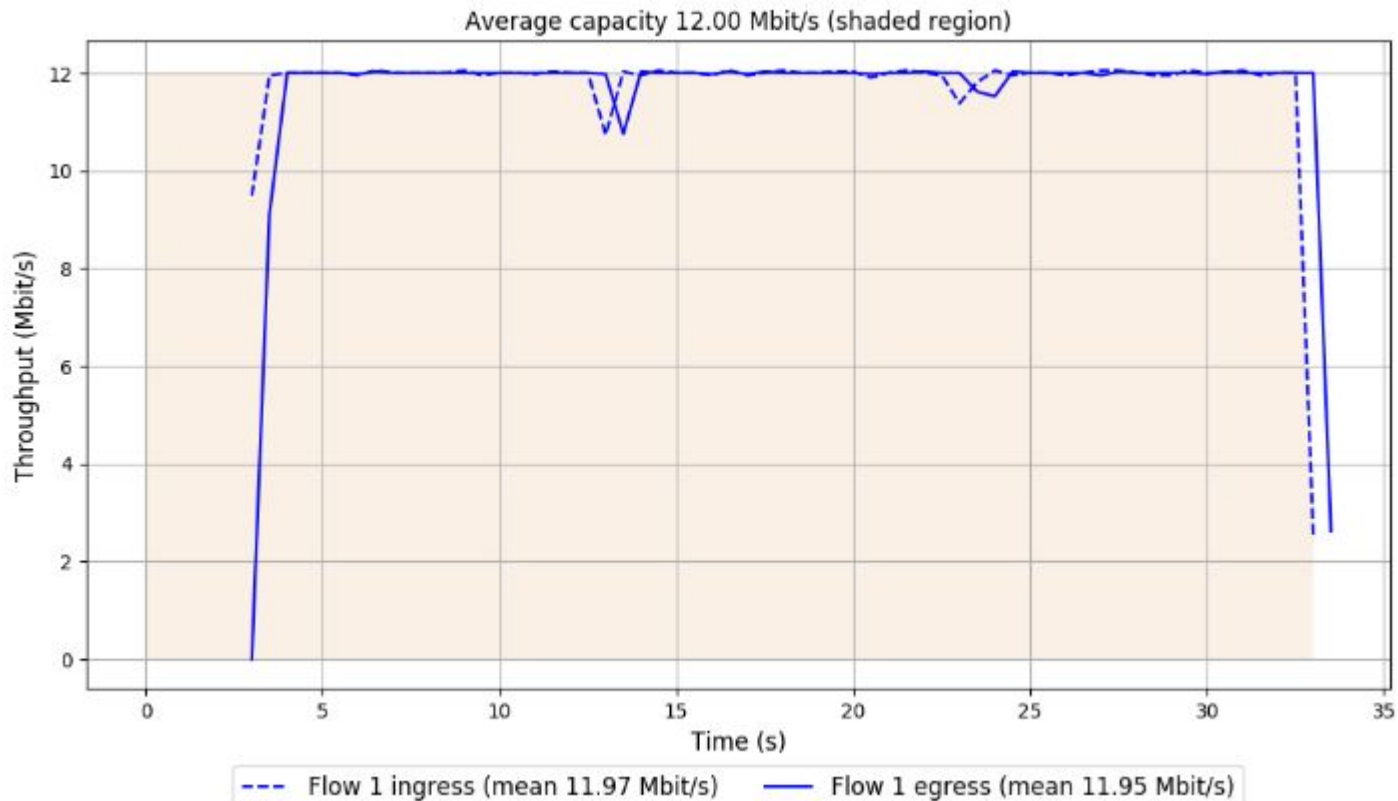
- Implement paper on Multi-Objective Congestion Control.
- Integrate MOCC with Orca and DeepCC (if time permits).
- Study performance of new architecture with different RL algorithms - DDPG, TD3, PPO, SAC, etc.

Pantheon

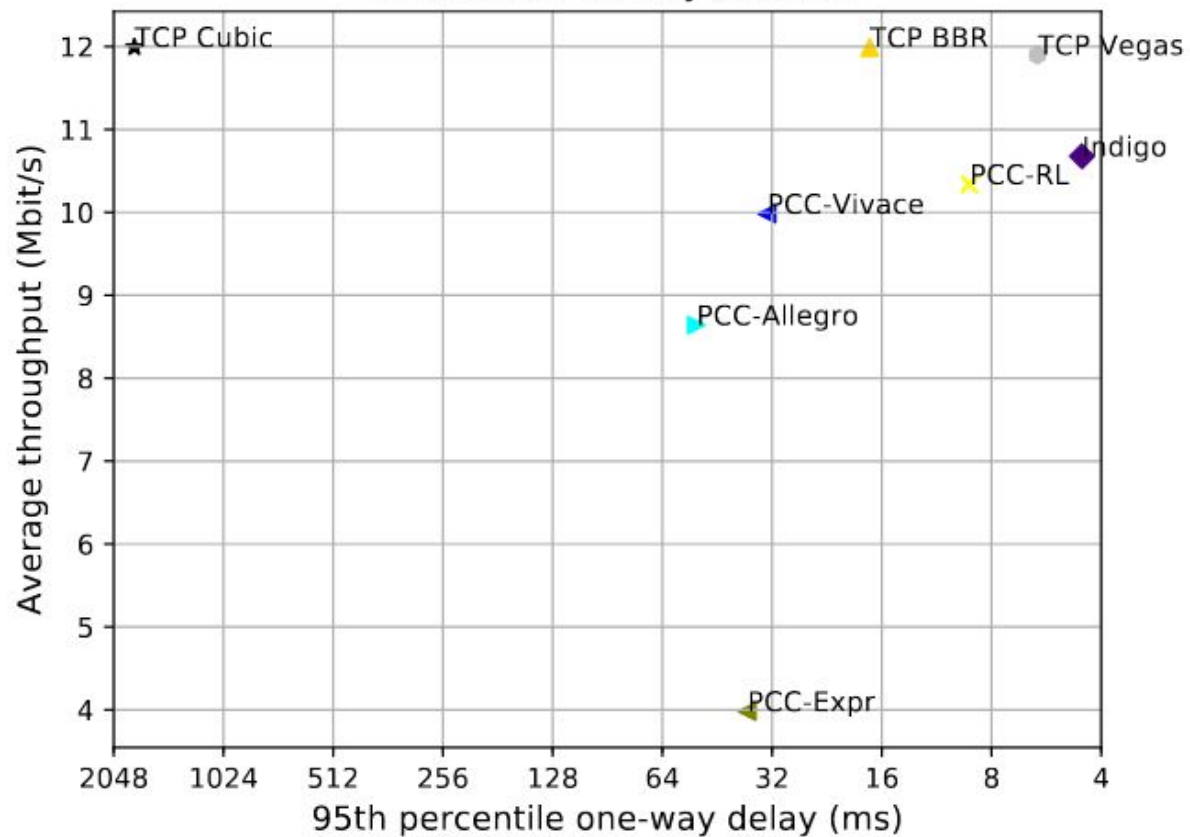
- Pantheon of Congestion Control: It is a community evaluation platform for academic research on congestion control that reduces the need to reinvent the wheel in the evaluations of new internet congestion-control algorithms.
- Pantheon provides:
 - 1) A collection of 17 working implementations of congestion-control schemes, all of them continuously verified to compile and run by a continuous integration system.
 - 2) A testbed of measurement nodes on wired and cellular networks.
 - 3) A collection of network emulators that can be used to test congestion-control schemes locally.
- The Pantheon performs several types of measurements on a roughly weekly basis. All measurements run a particular congestion-control scheme between two endpoints, measuring the departure time of each IP datagram (at the sender) and the arrival time of the same IP datagram (at the receiver), if it arrives. These raw logs are available for each measurement. For each scheme, it also calculates and plots aggregate statistics, e.g., the throughput, one-way delay (95th percentile), loss rate, etc.
- To run a new scheme, submit a pull request, the Travis-CI system will automatically verify that the scheme compiles and runs in emulation.

Test 1

12 Mbps trace

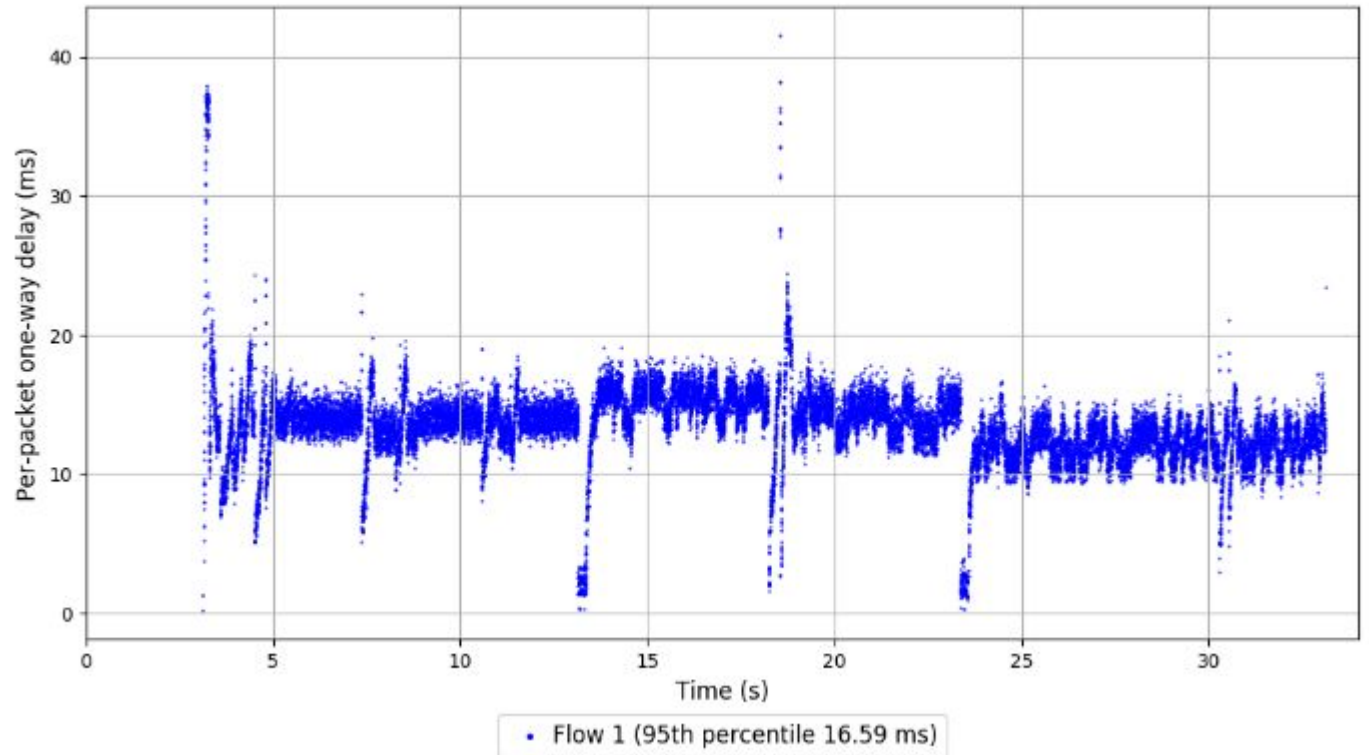


local test in mahimahi, 1 run of 30s each per scheme
(mean of all runs by scheme)

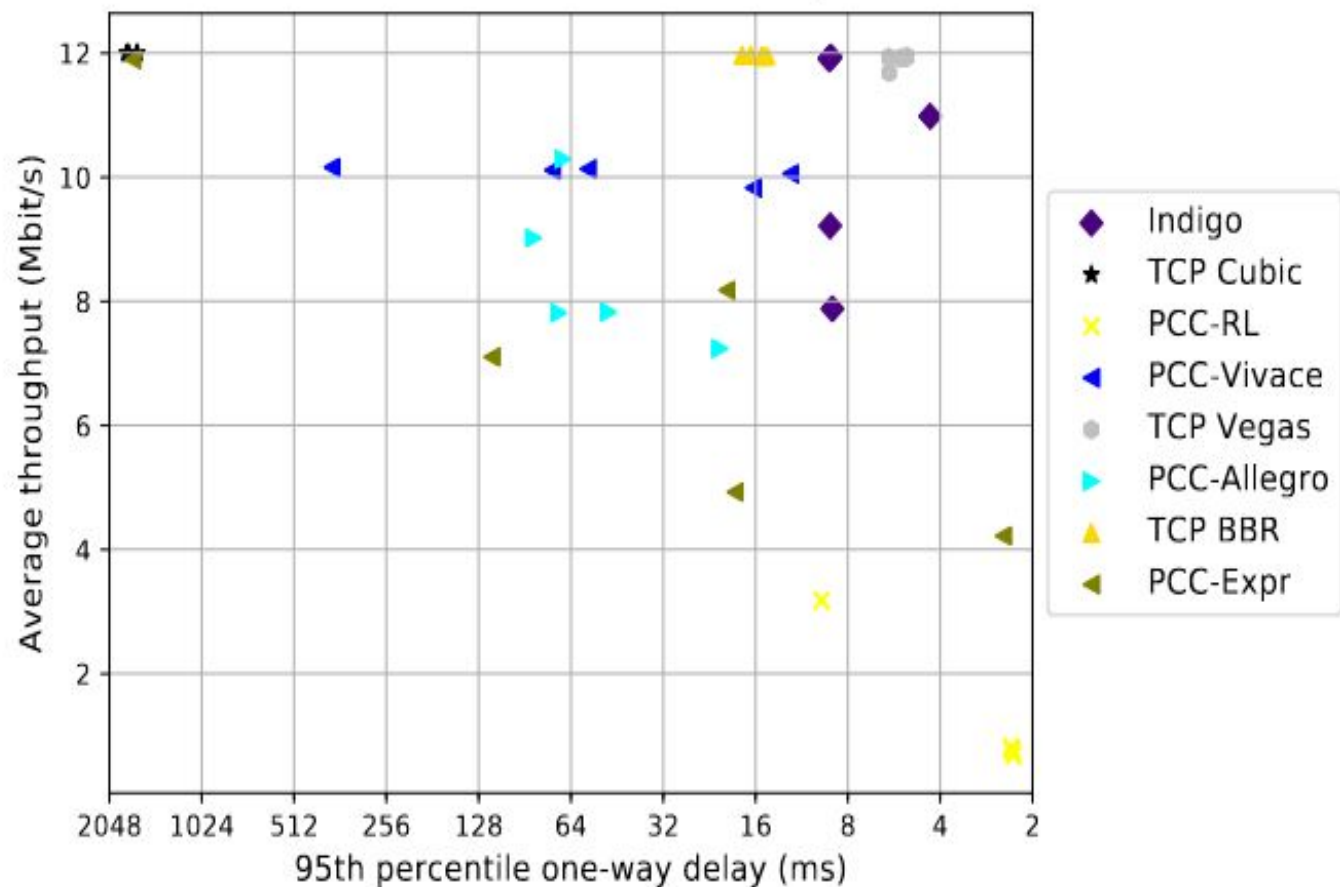


Test 2

Multiple runs

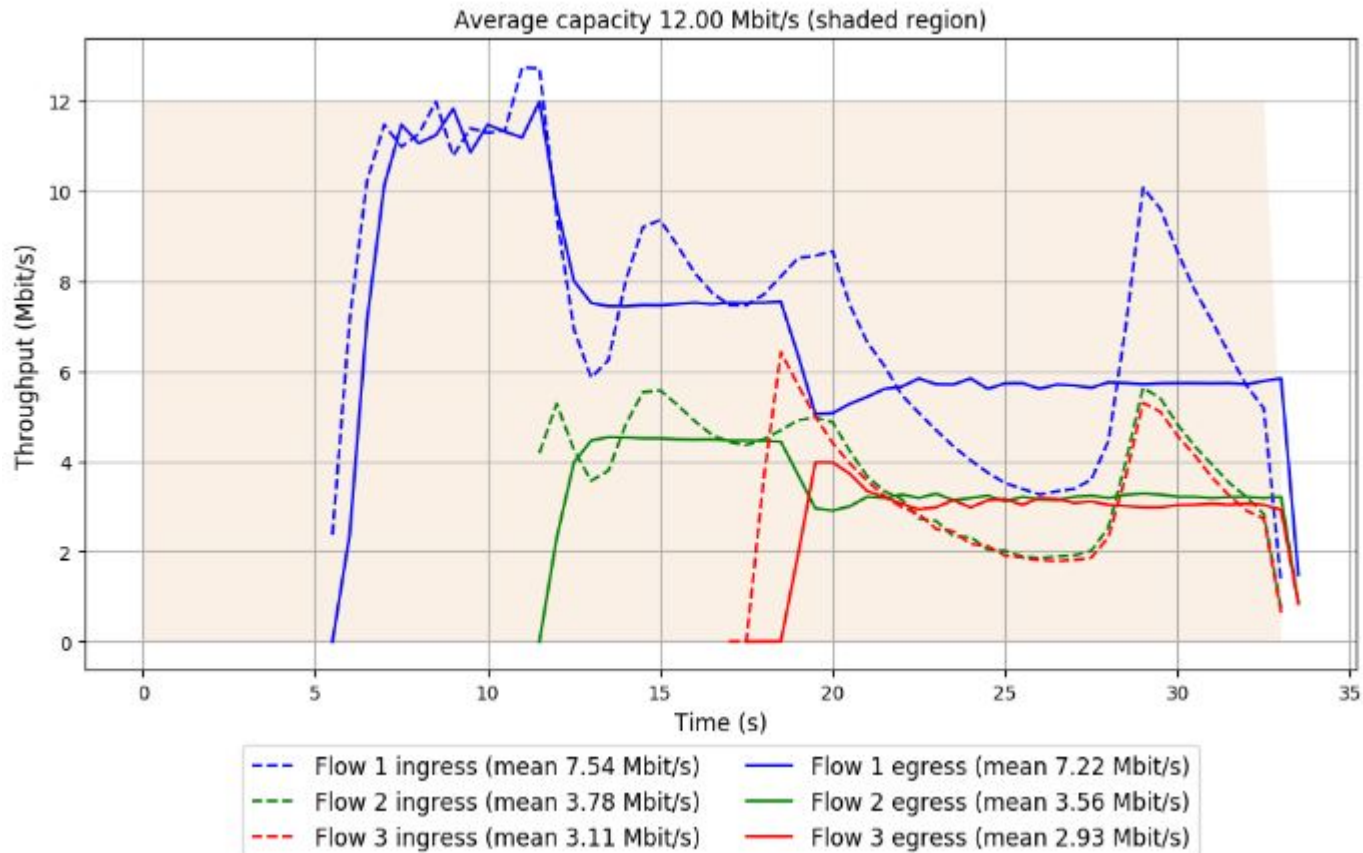


local test in mahimahi, 5 runs of 30s each per scheme

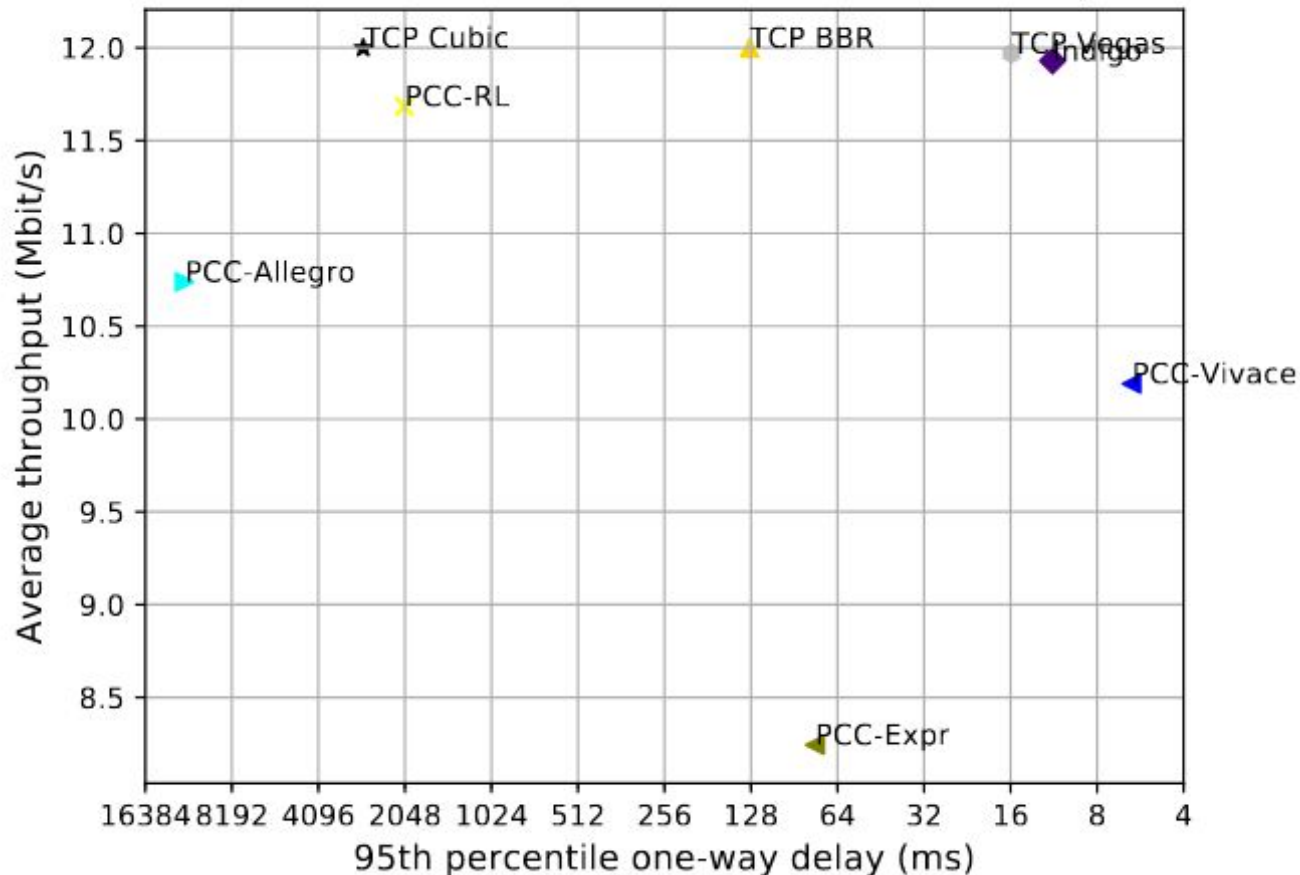


Test 3

Multiple flows

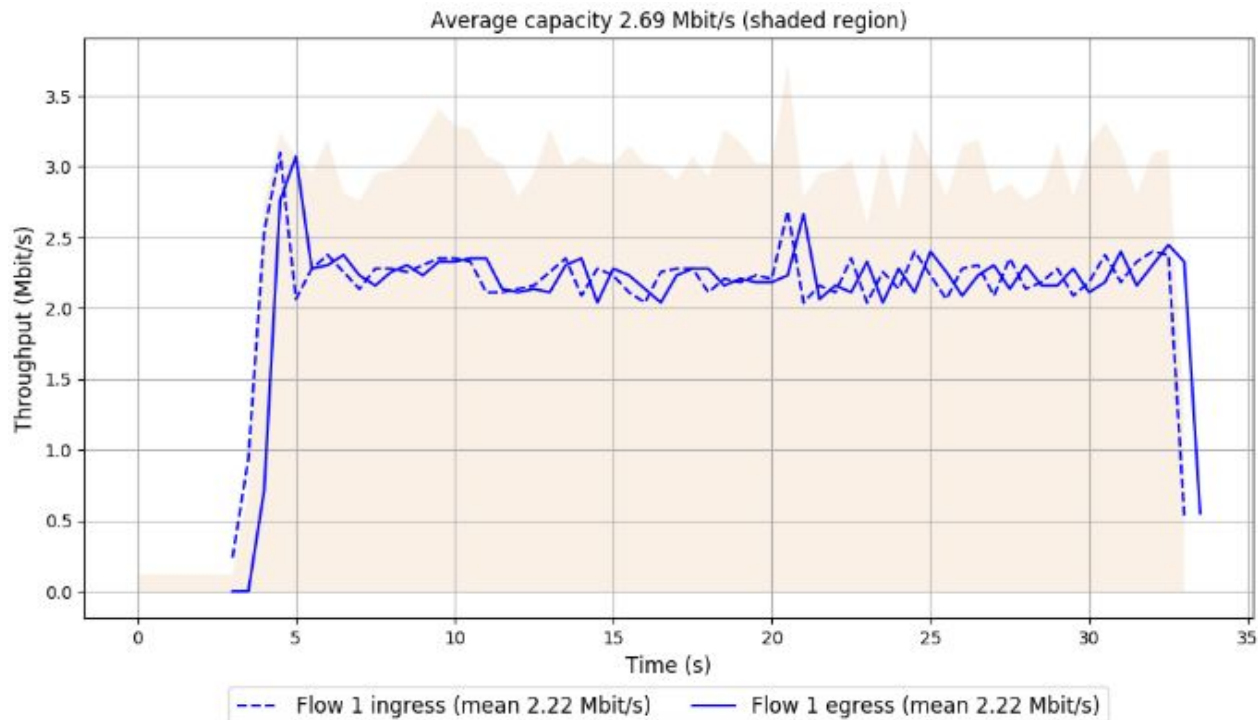


local test in mahimahi, 1 run of 30s each per scheme
3 flows with 5s interval between flows (mean of all runs by scheme)

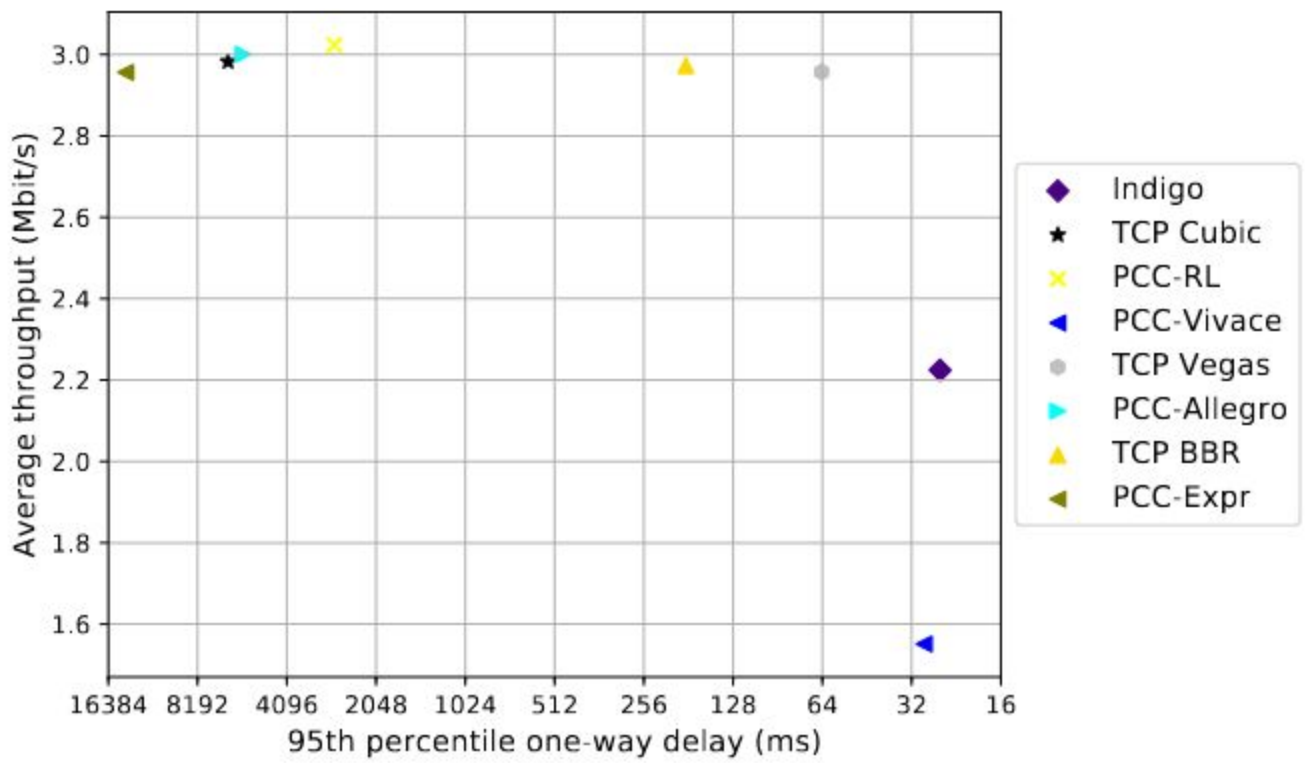


Test 4

Poisson distributed trace

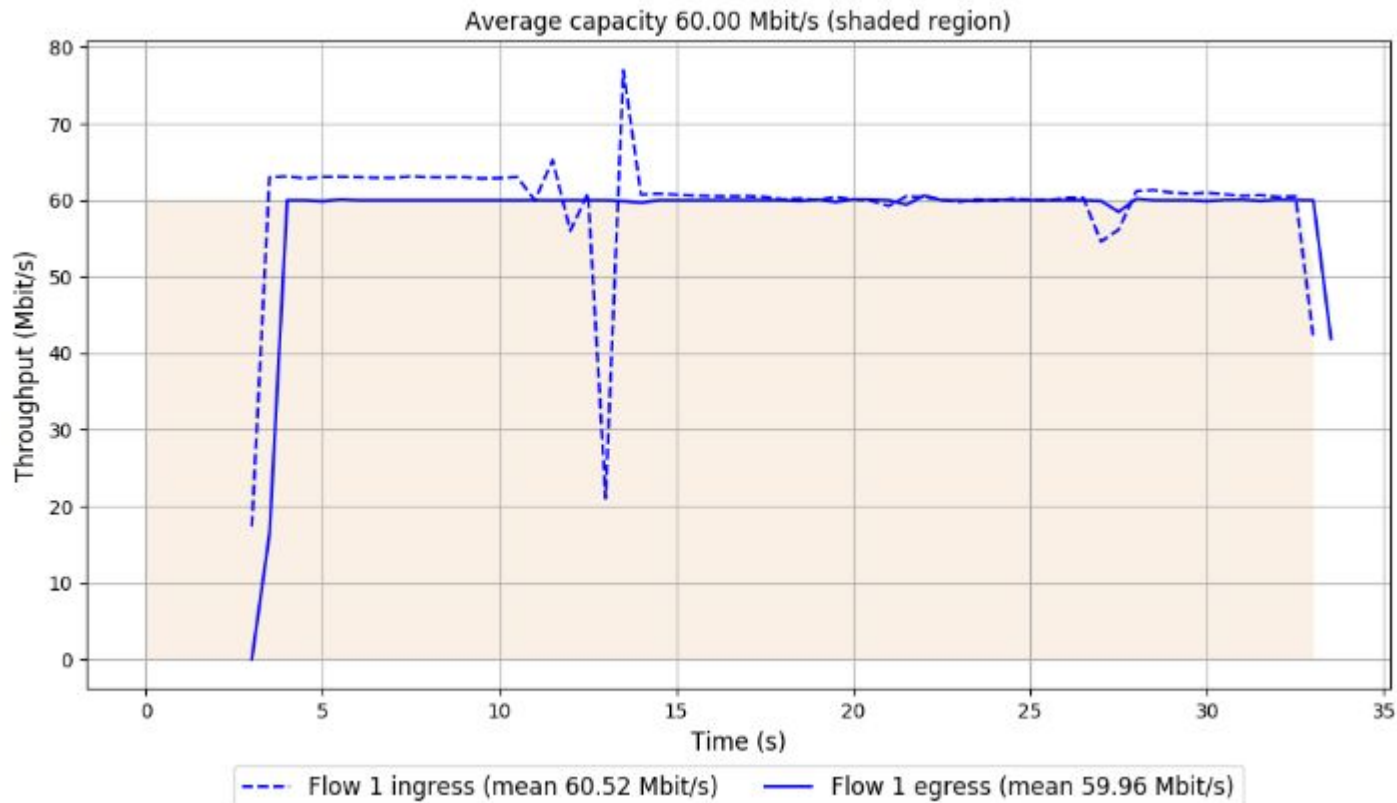


local test in mahimahi, 1 run of 30s each per scheme

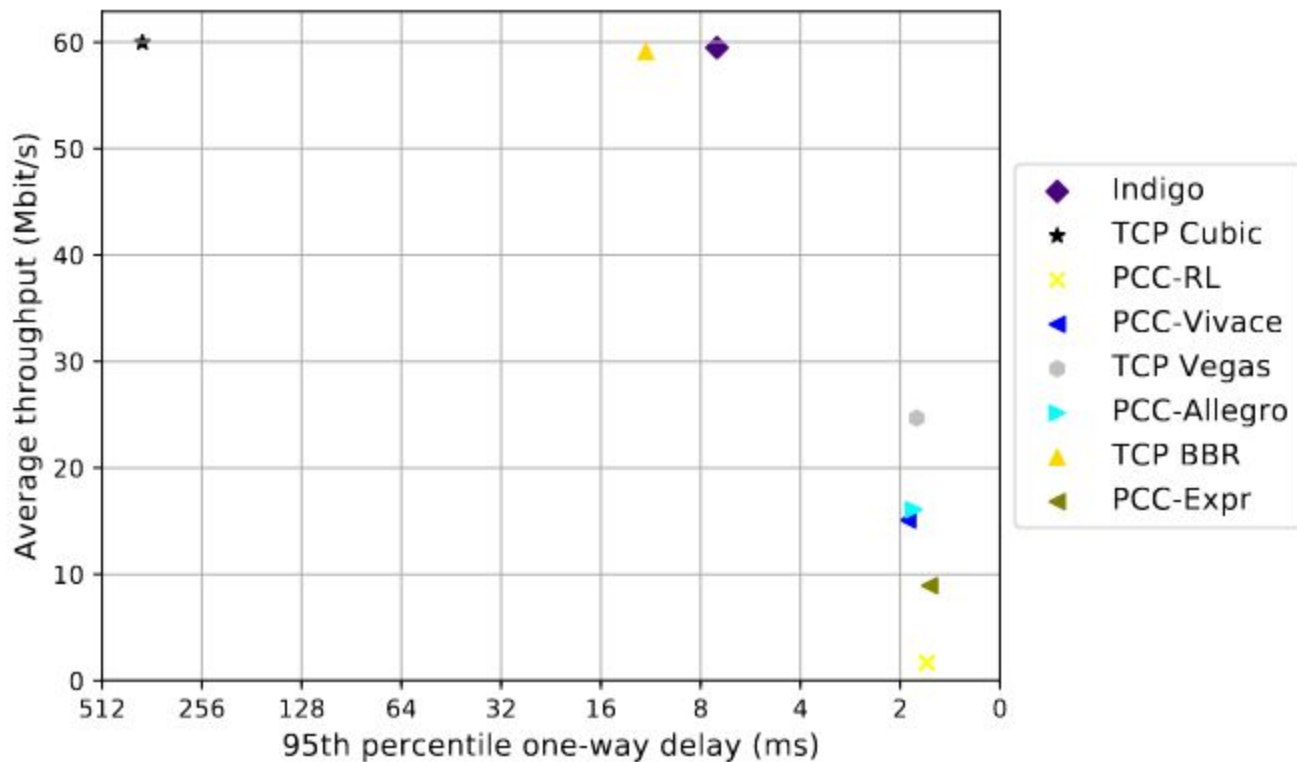


Test 5

60 Mbps trace

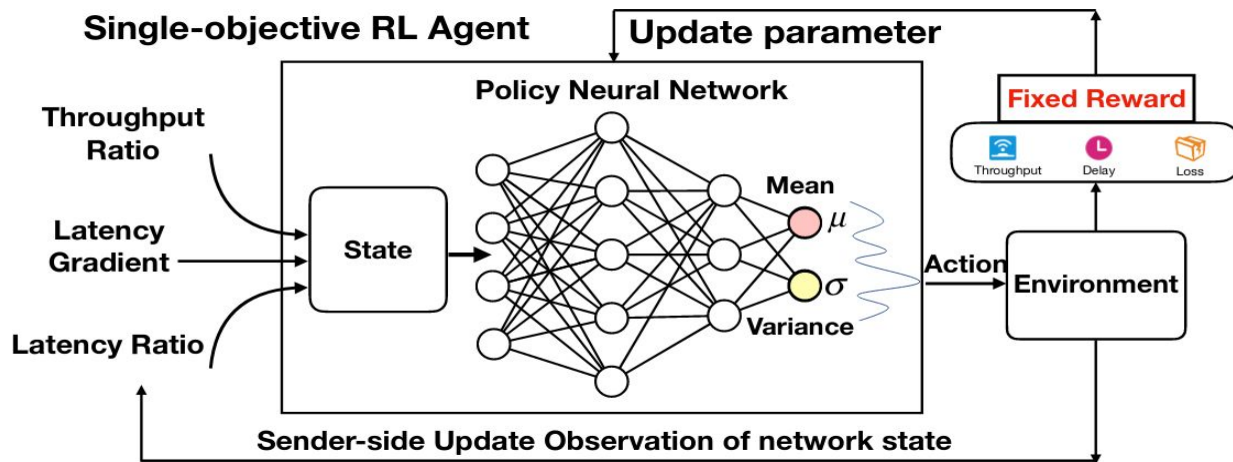


local test in mahimahi, 1 run of 30s each per scheme



Aurora

- RL Algorithm: Proximal Policy Optimization (PPO1)
- Architecture (of NN): [32, 16]
- History Length: 10
- Features: ['sent latency inflation', 'latency ratio', 'send ratio']
- Gamma: 0.99

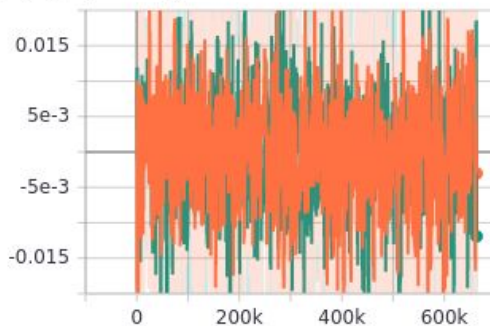


Aurora: Implementation Details

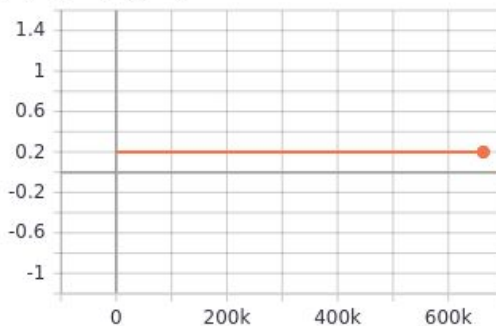
- States: Fixed-length history of statistics vector - latency gradient, latency ratio, sending ratio
- Actions: Adjust sending rate every monitoring interval
- Reward: $10 * \textit{throughput} - 1000 * \textit{latency} - 2000 * \textit{loss}$
- RL Algorithm: Proximal Policy Optimization. The PPO algorithm combines ideas from A2C (having multiple workers) and TRPO (it uses a trust region to improve the actor). The main idea is that after an update, the main policy should not be too far from the old policy. For that, PPO uses clipping to avoid too large updates. It empirically performs at least as close to TRPO.

PPO [64, 32]: Training Statistics

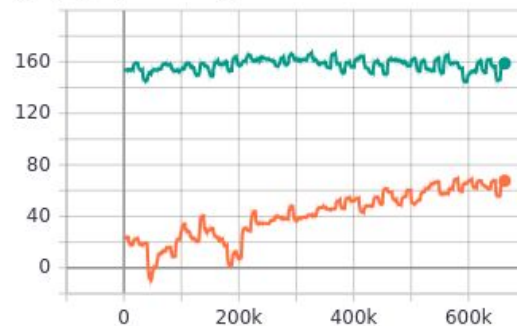
advantage
tag: input_info/advantage



clip_range
tag: input_info/clip_range



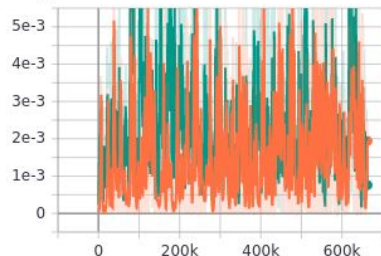
discounted_rewards
tag: input_info/discounted_rewards



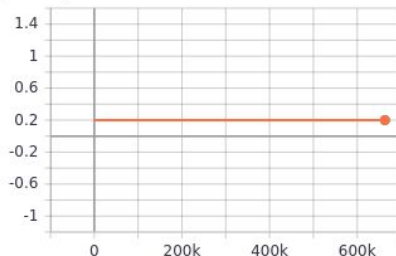
Orange graph - First checkpoint, Green graph - Last checkpoint

PPO [64, 32]: Training Statistics

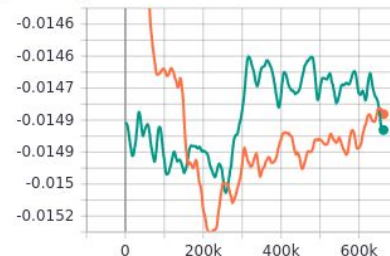
approximate_kullback-leibler
tag: loss/approximate_kullback-leibler



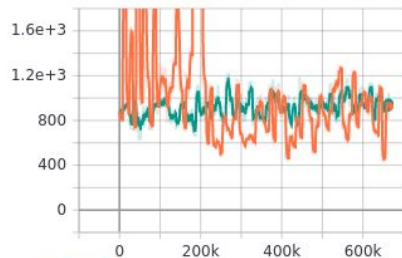
clip_factor
tag: loss/clip_factor



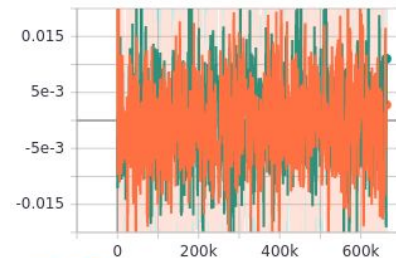
entropy_loss
tag: loss/entropy_loss



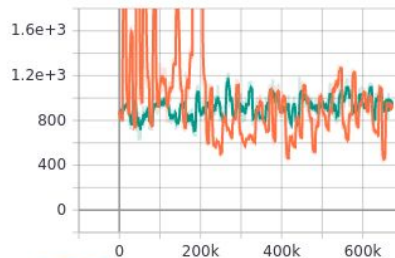
loss
tag: loss/loss



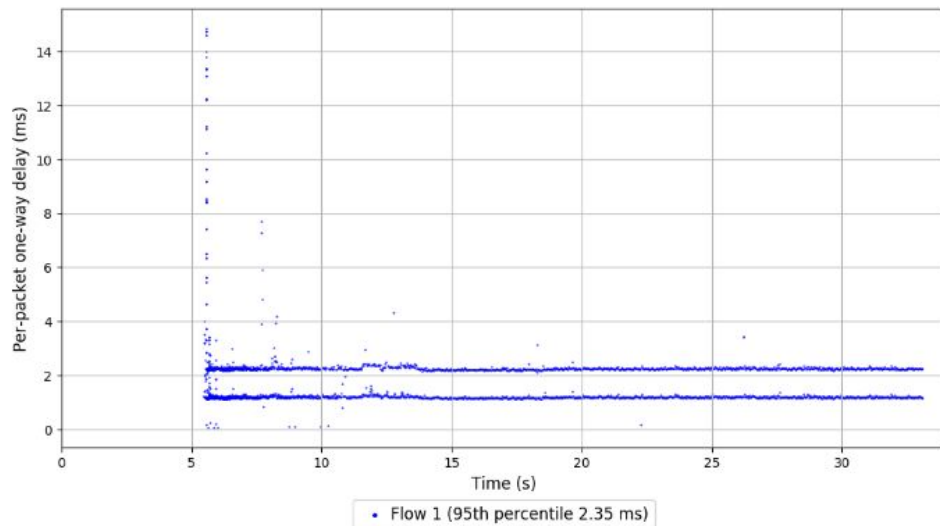
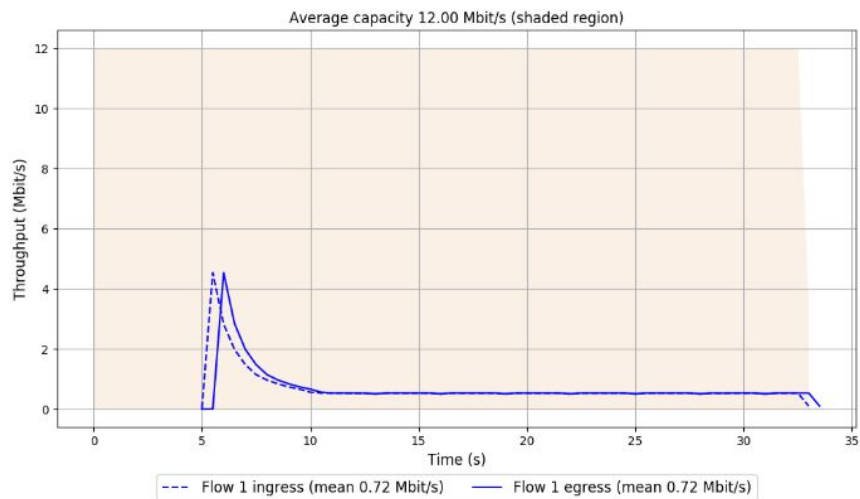
policy_gradient_loss
tag: loss/policy_gradient_loss



value_function_loss
tag: loss/value_function_loss

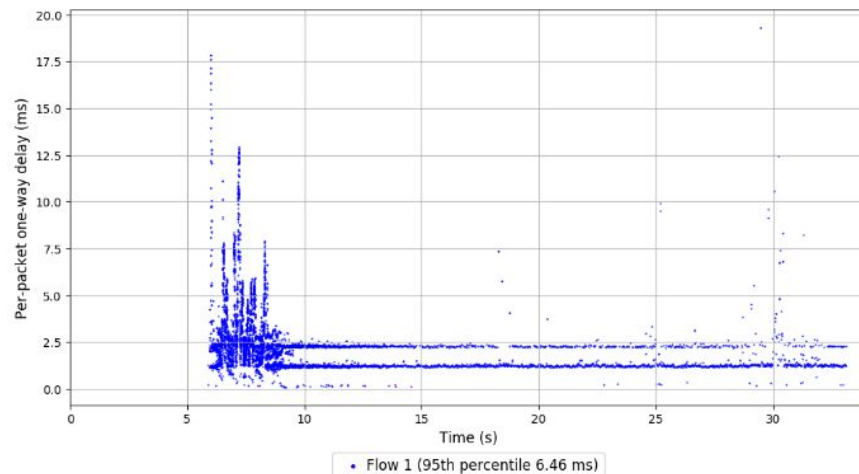
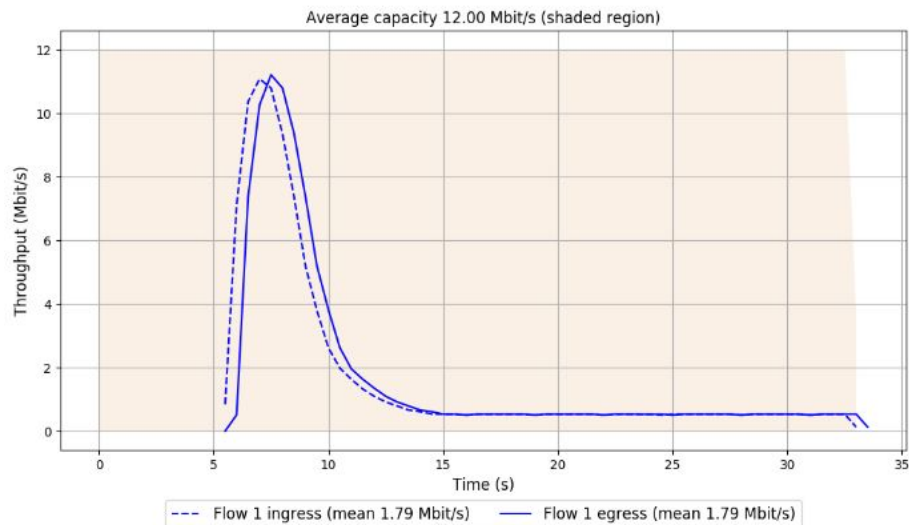


Test 1: PPO [64, 32]



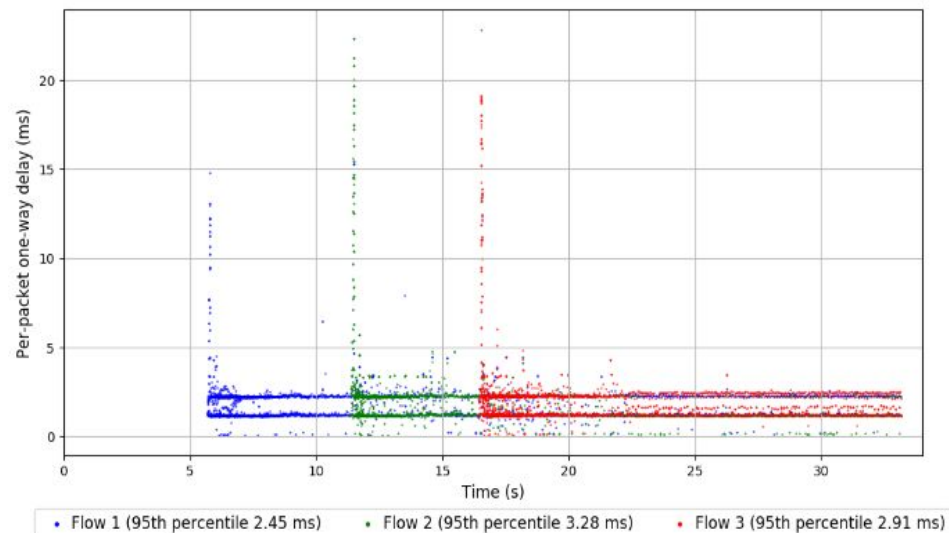
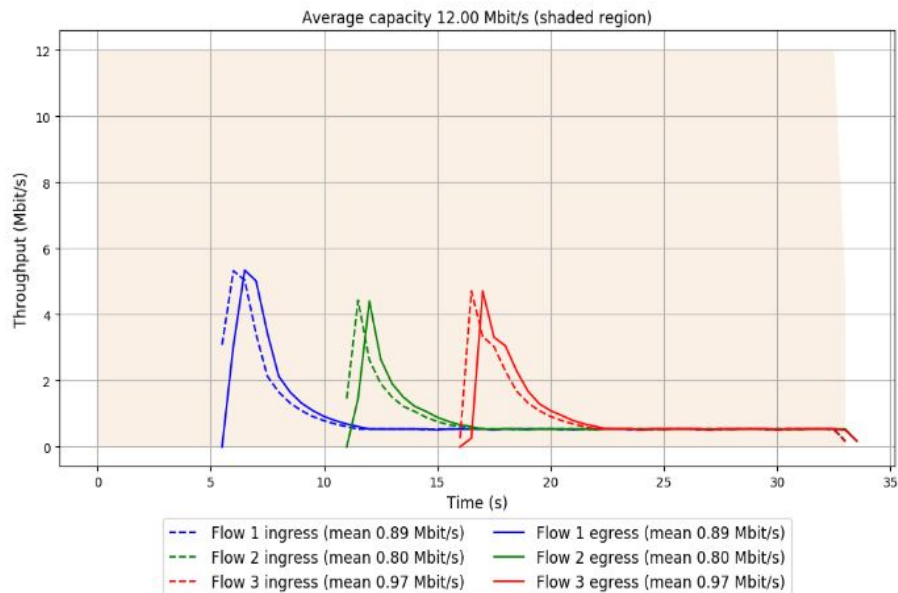
scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	0.72	2.35	0.00

Test 2: PPO [64, 32]



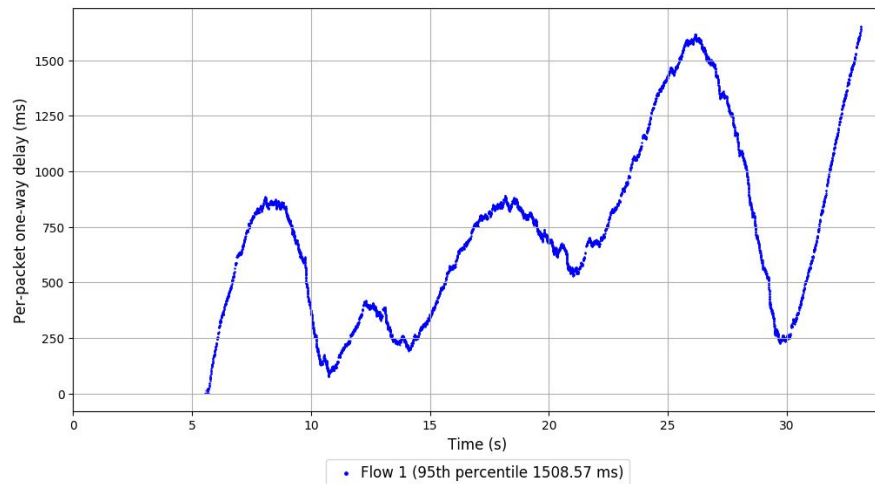
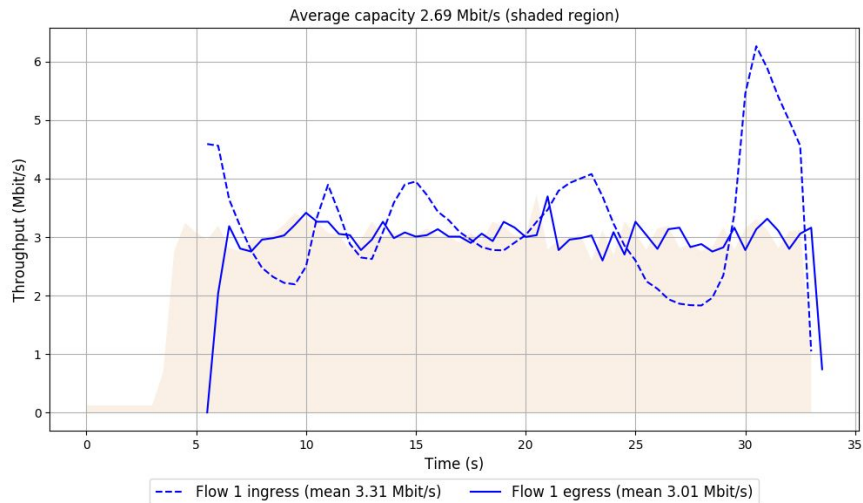
scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	5	0.98	3.25	0.00

Test 3: PPO [64, 32]



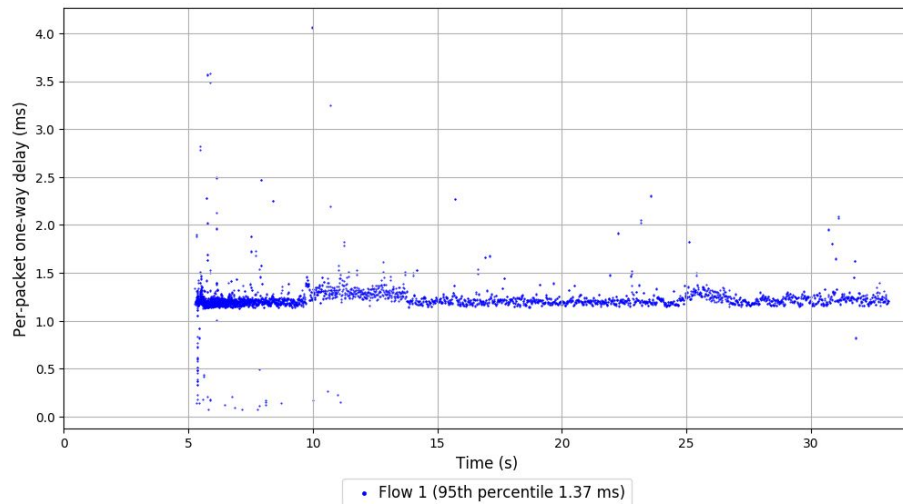
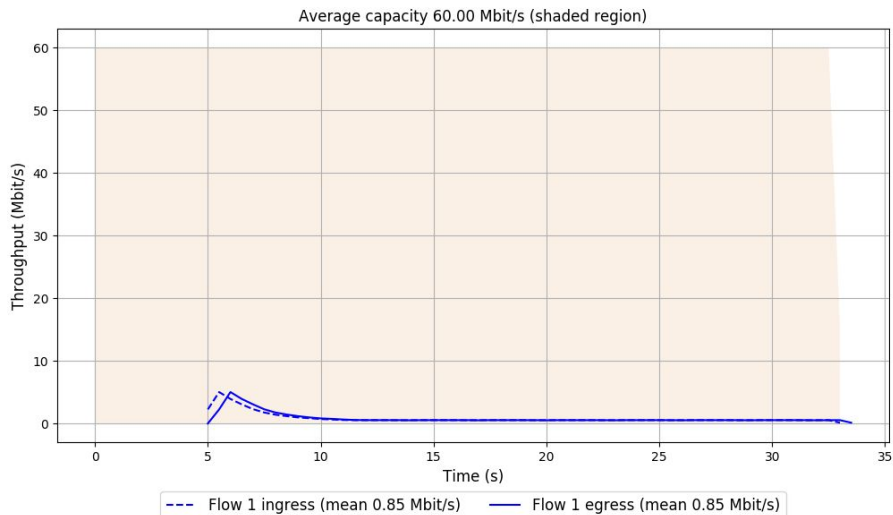
scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	0.89	0.80	0.97	2.45	3.28	2.91	0.00	0.00	0.00

Test 4: PPO [64, 32]



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	3.01	1508.57	8.92

Test 5: PPO [64,32]



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	0.85	1.37	0.00

Aurora: Network Architecture Tests

RL Algorithm: Proximal Policy Optimization (PPO1)

Architecture (of NN): [32, 16] → [**128, 8**]

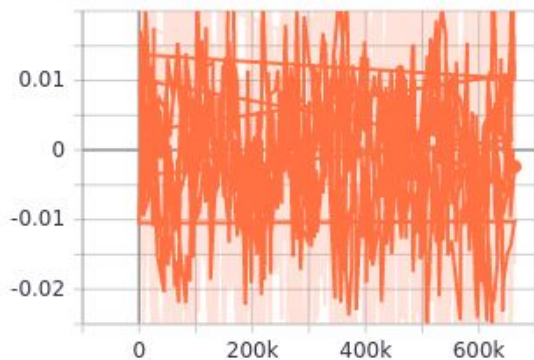
History Length: 10

Features: ['sent latency inflation', 'latency ratio', 'send ratio']

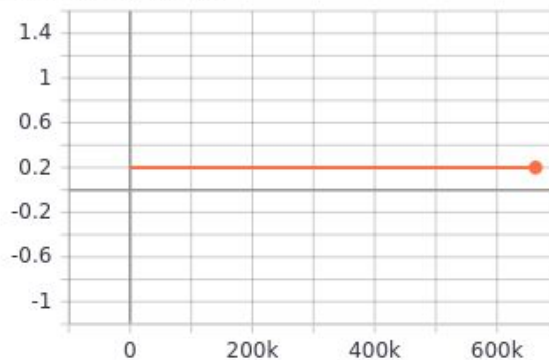
Gamma: 0.99

PPO [128, 8]: Training Statistics

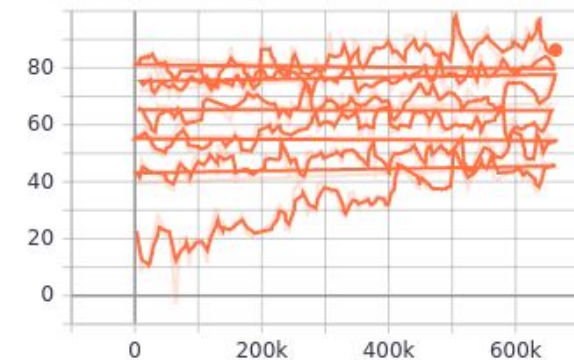
advantage
tag: input_info/advantage



clip_range
tag: input_info/clip_range



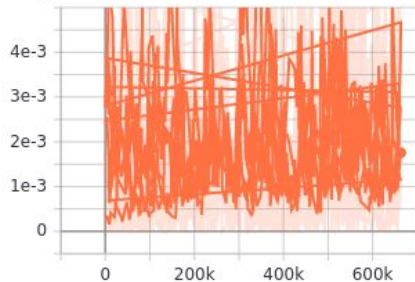
discounted_rewards
tag: input_info/discounted_rewards



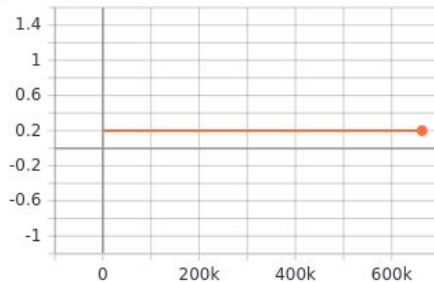
All six training checkpoints - discounted rewards increase with each checkpoint

PPO [128, 8]: Training Statistics

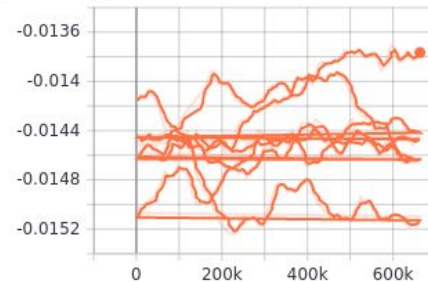
approximate_kullback-leibler
tag: loss/approximate_kullback-leibler



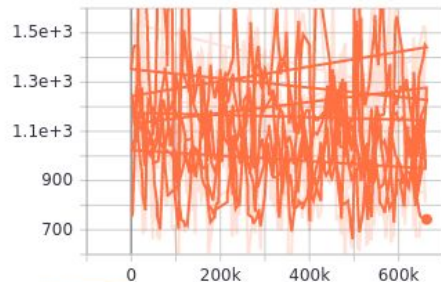
clip_factor
tag: loss/clip_factor



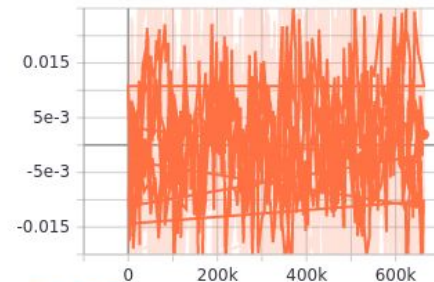
entropy_loss
tag: loss/entropy_loss



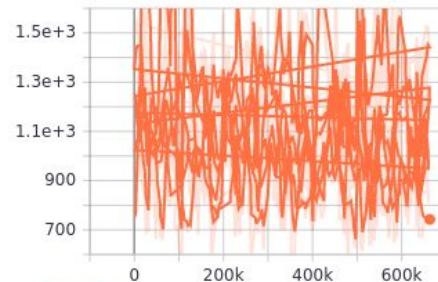
loss
tag: loss/loss



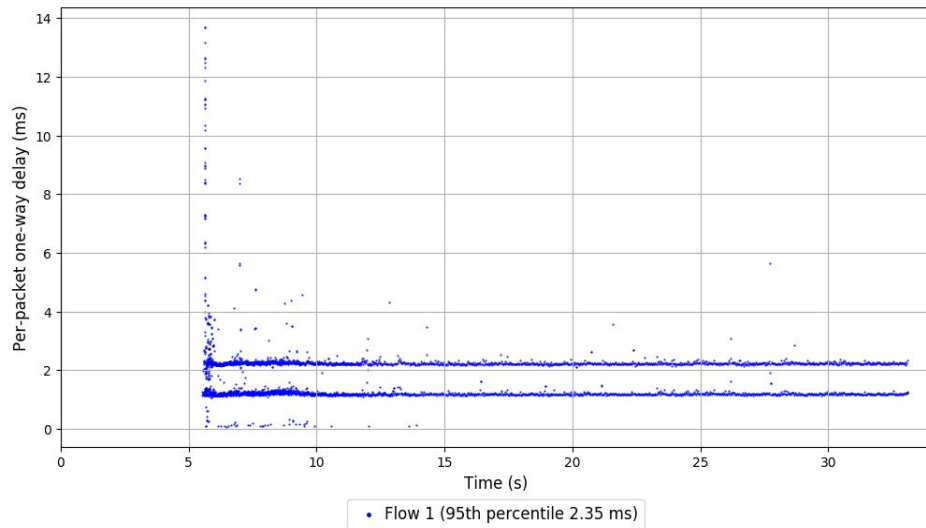
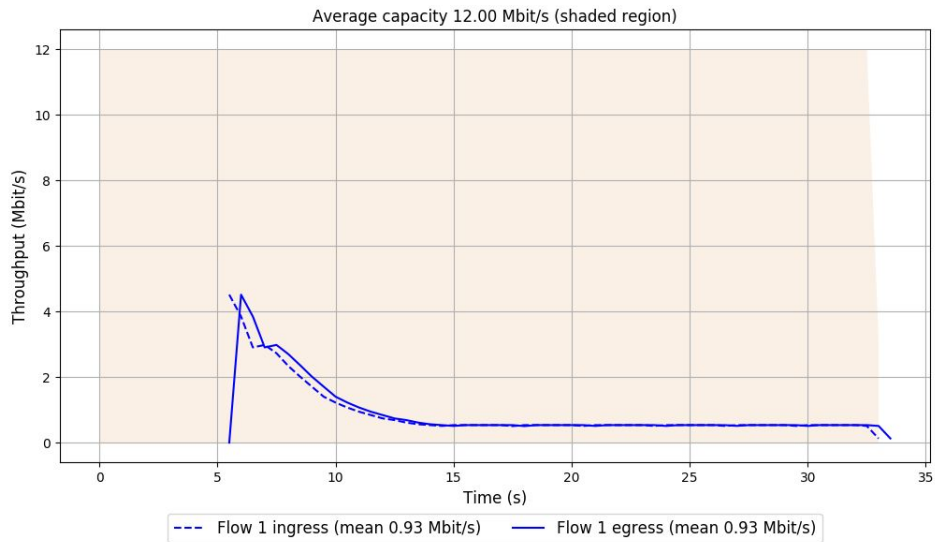
policy_gradient_loss
tag: loss/policy_gradient_loss



value_function_loss
tag: loss/value_function_loss

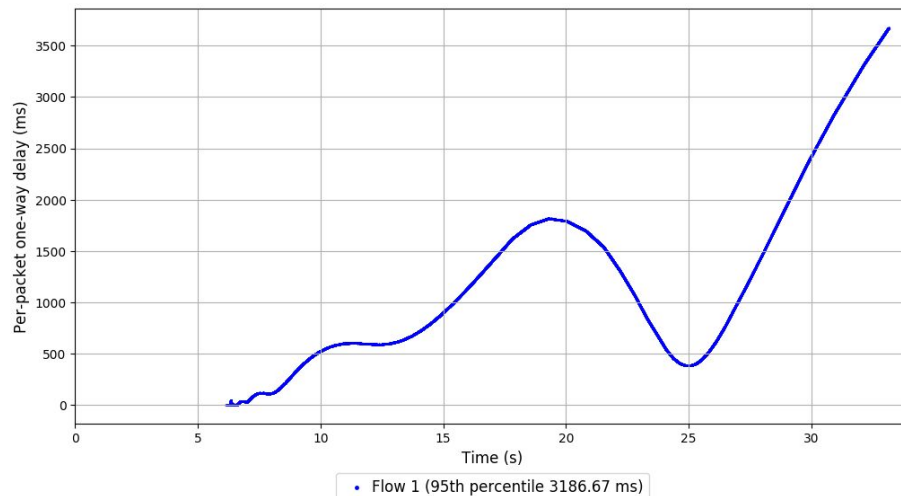
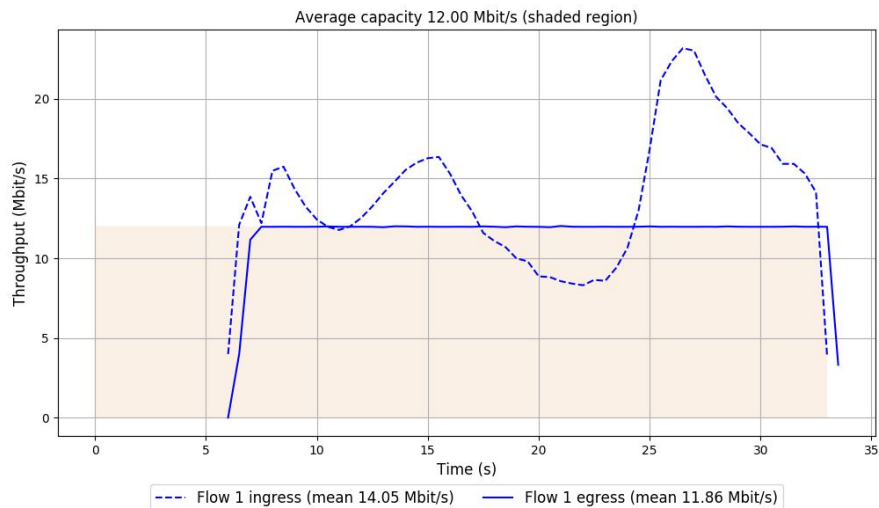


Test 1: PPO [128, 8]



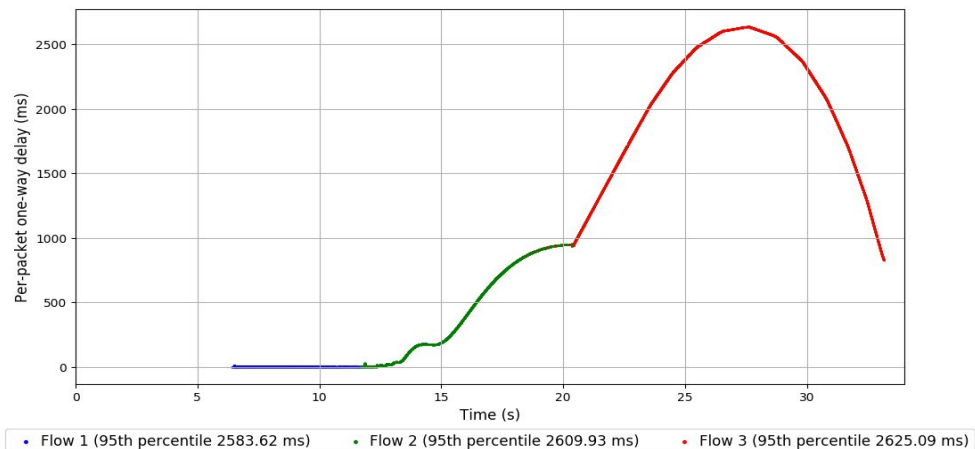
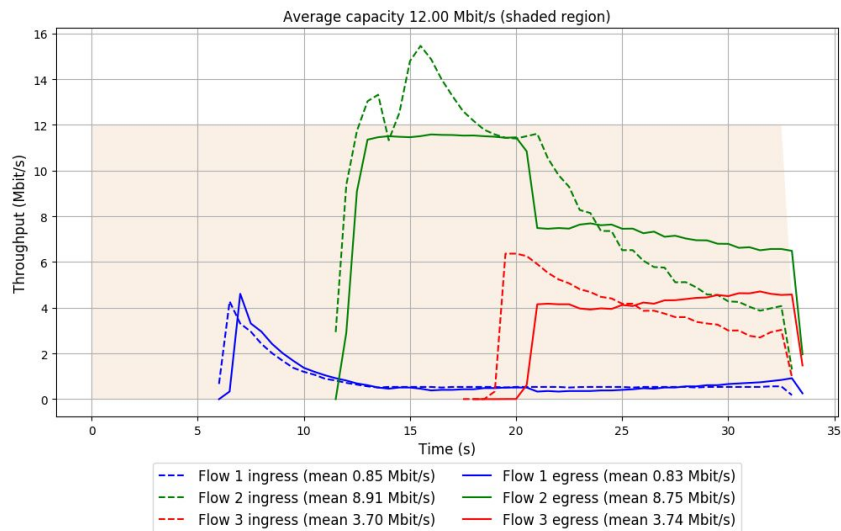
scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	1	0.93	2.35	0.00

Test 2: PPO [128, 8]



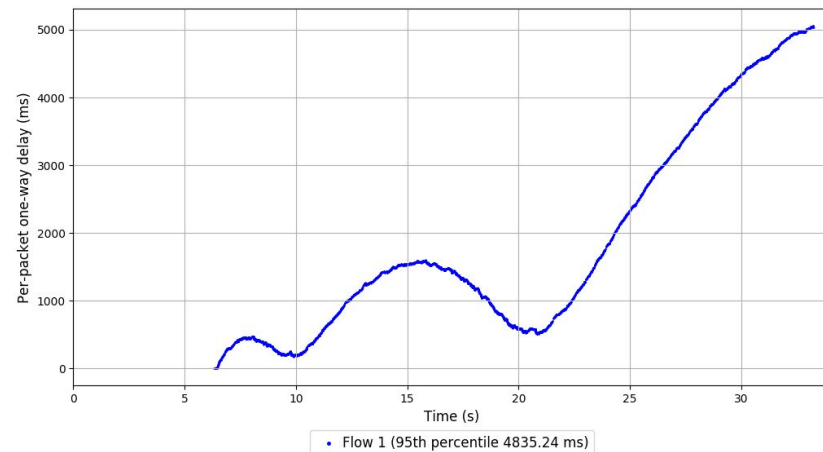
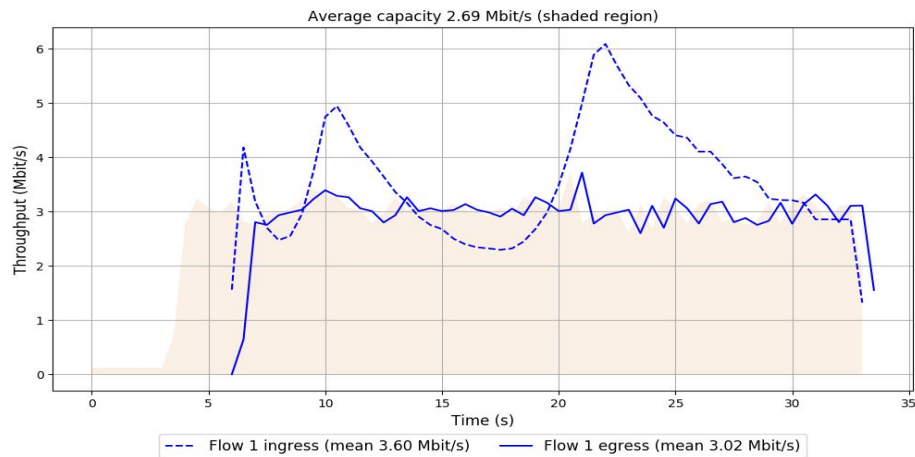
scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	5	9.63	2880.08	11.33

Test 3: PPO [128, 8]



scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	0.83	8.75	3.74	2583.62	2609.93	2625.09	2.07	1.77	4.35

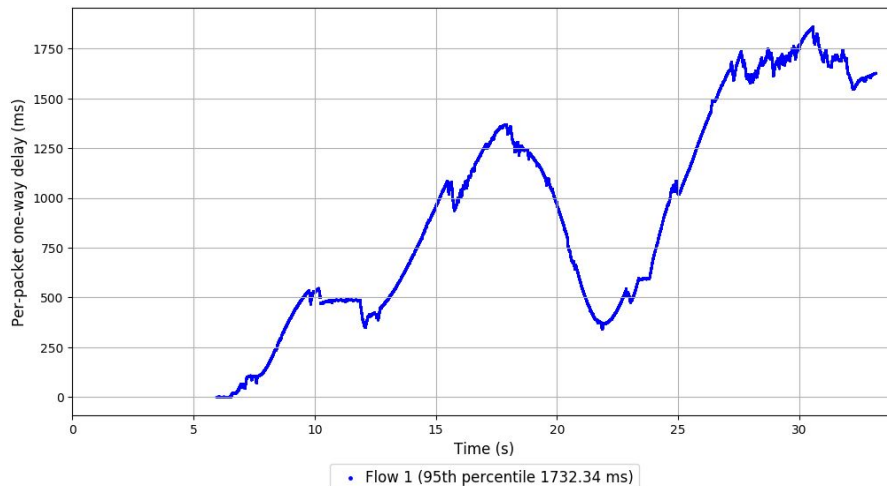
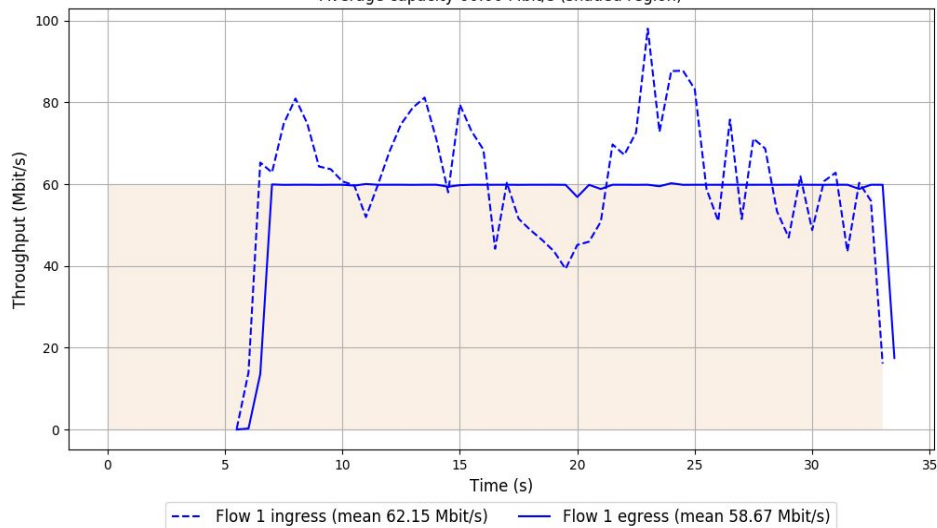
Test 4: PPO [128, 8]



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	3.02	4835.24	16.13

Test 5: PPO [128,8]

Average capacity 60.00 Mbit/s (shaded region)



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	58.67	1732.34	5.60

Aurora: Network Architecture Tests

RL Algorithm: Proximal Policy Optimization (PPO1)

Architecture (of NN): [32, 16] → [**64, 32, 64**]

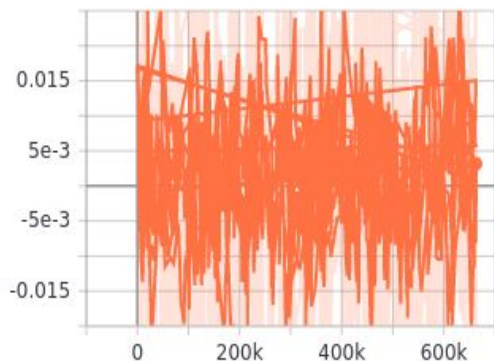
History Length: 10

Features: ['sent latency inflation', 'latency ratio', 'send ratio']

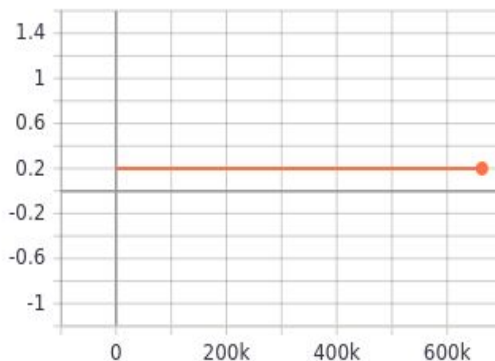
Gamma: 0.99

PPO [64, 32, 64]: Training Statistics

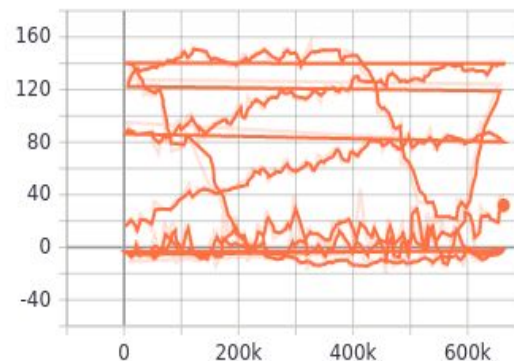
advantage
tag: input_info/advantage



clip_range
tag: input_info/clip_range



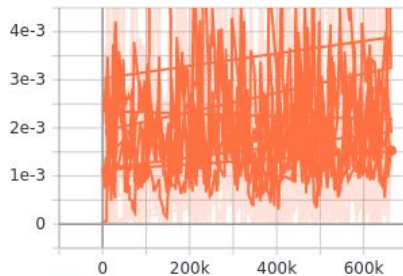
discounted_rewards
tag: input_info/discounted_rewards



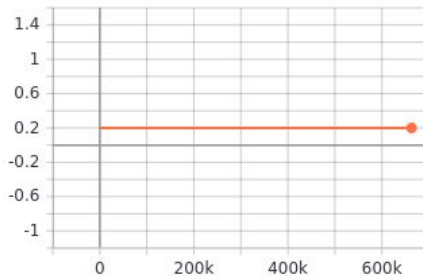
All six training checkpoints - discounted rewards increase with each checkpoint

PPO [64, 32, 64]: Training Statistics

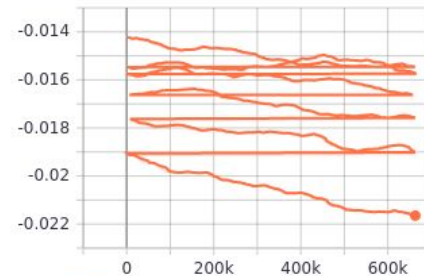
approximate_kullback-leibler
tag: loss/approximate_kullback-leibler



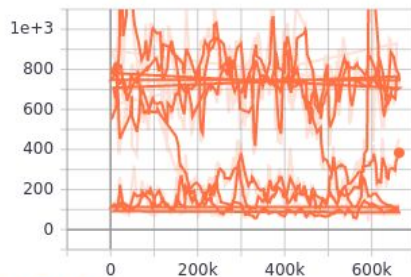
clip_factor
tag: loss/clip_factor



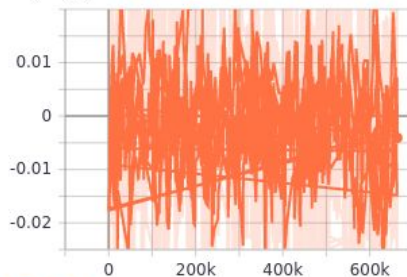
entropy_loss
tag: loss/entropy_loss



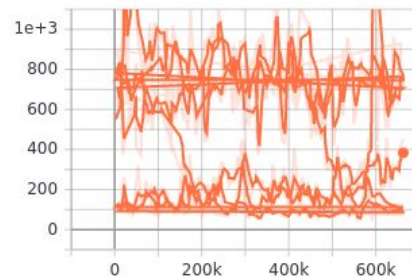
loss
tag: loss/loss



policy_gradient_loss
tag: loss/policy_gradient_loss

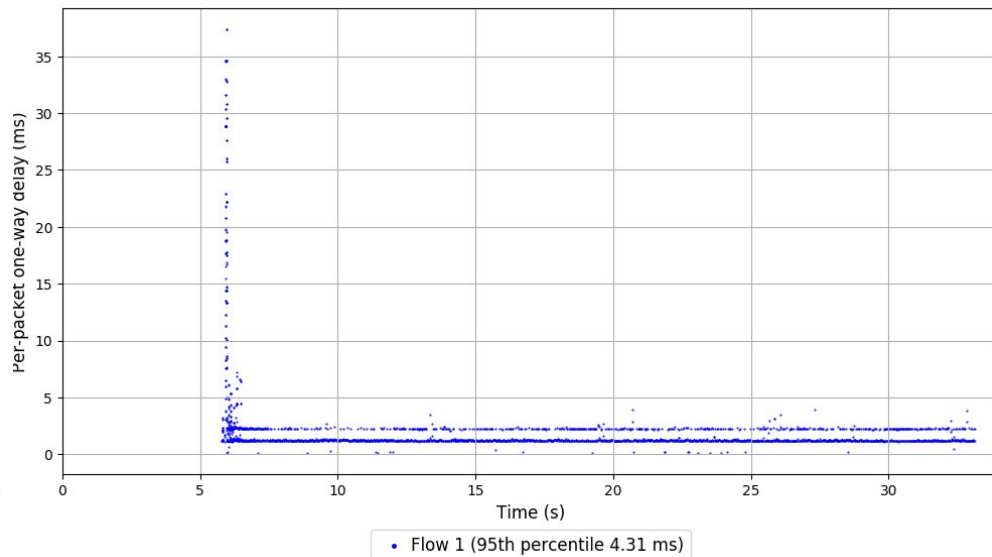
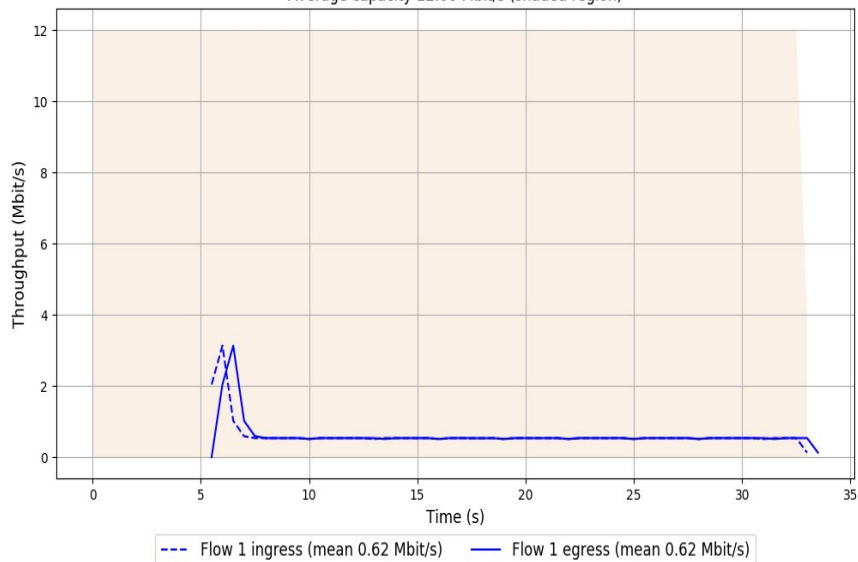


value_function_loss
tag: loss/value_function_loss



Test 1: PPO [64, 32, 64]

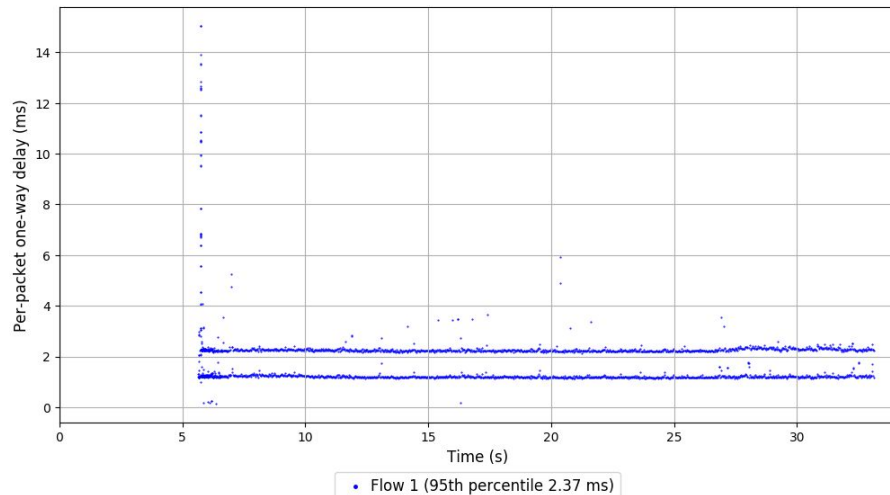
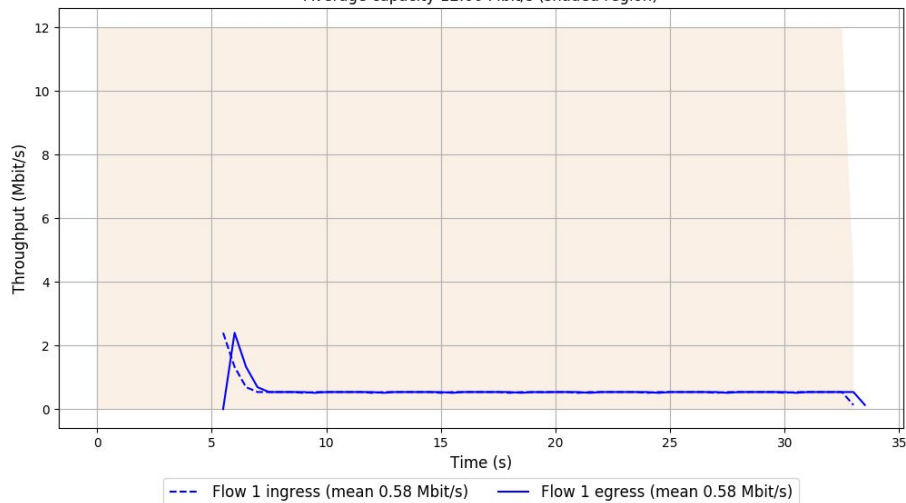
Average capacity 12.00 Mbit/s (shaded region)



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	0.62	4.31	0.00

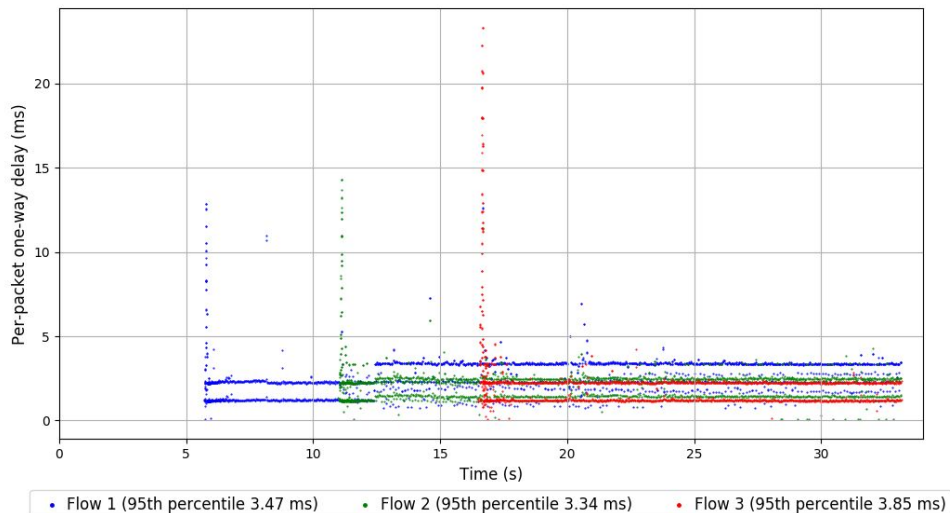
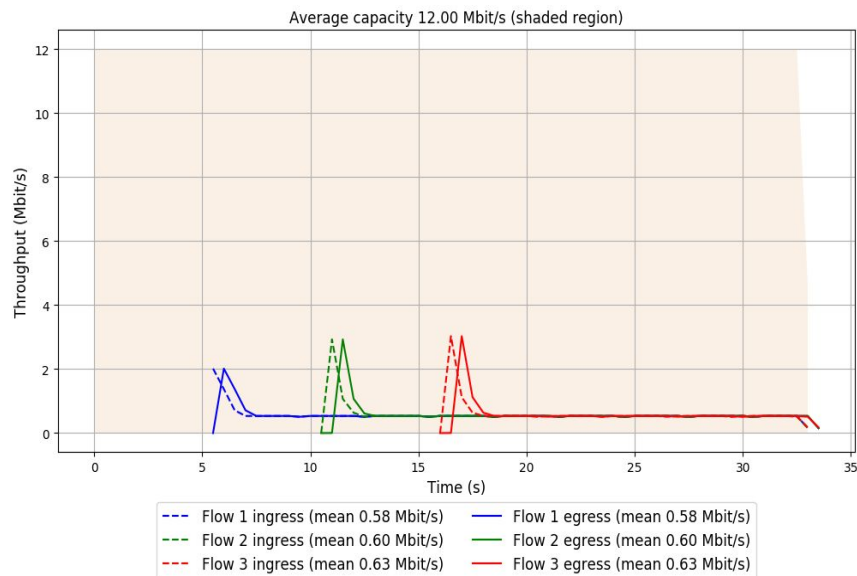
Test 2: PPO [64, 32, 64]

Average capacity 12.00 Mbit/s (shaded region)



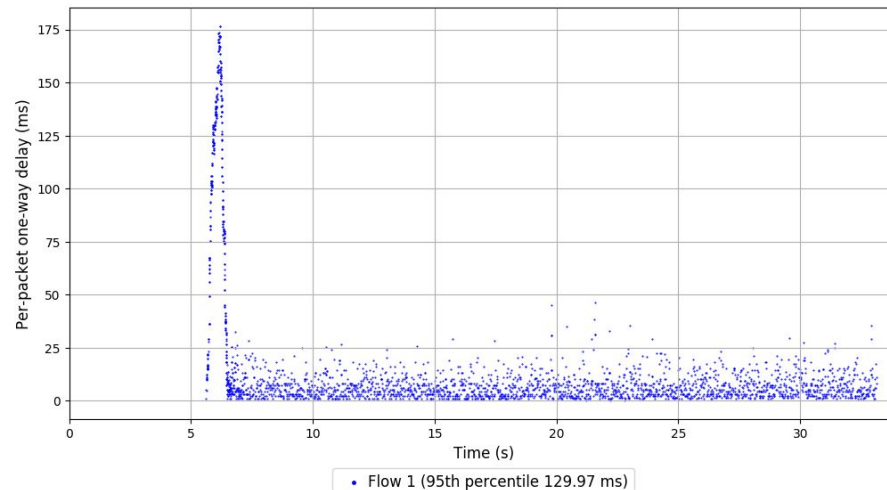
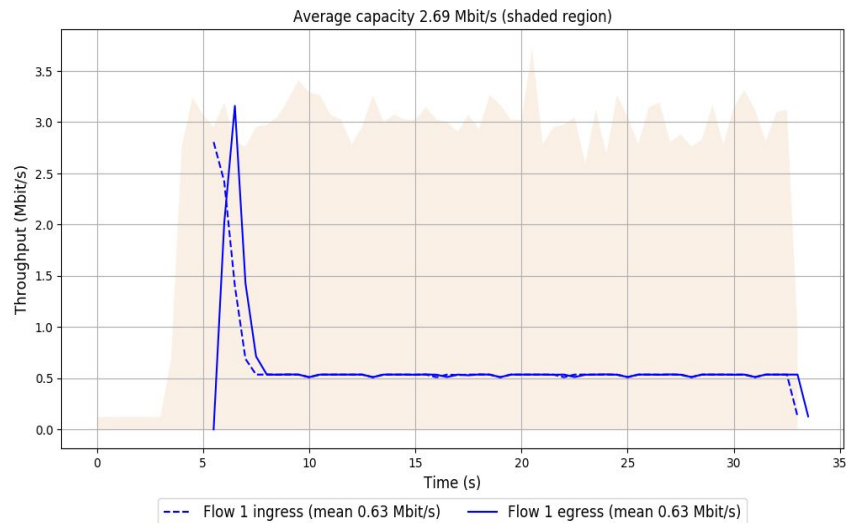
scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	5	0.59	2.46	0.00

Test 3: PPO [64, 32, 64]

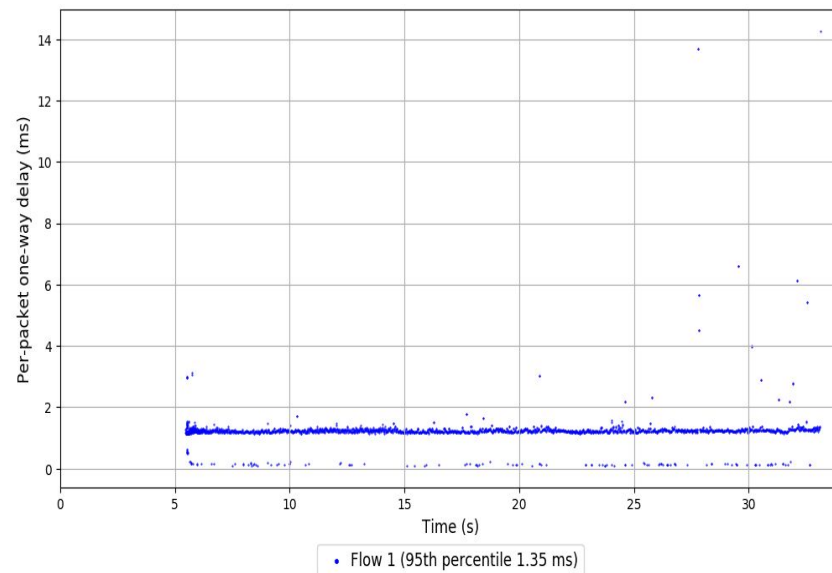
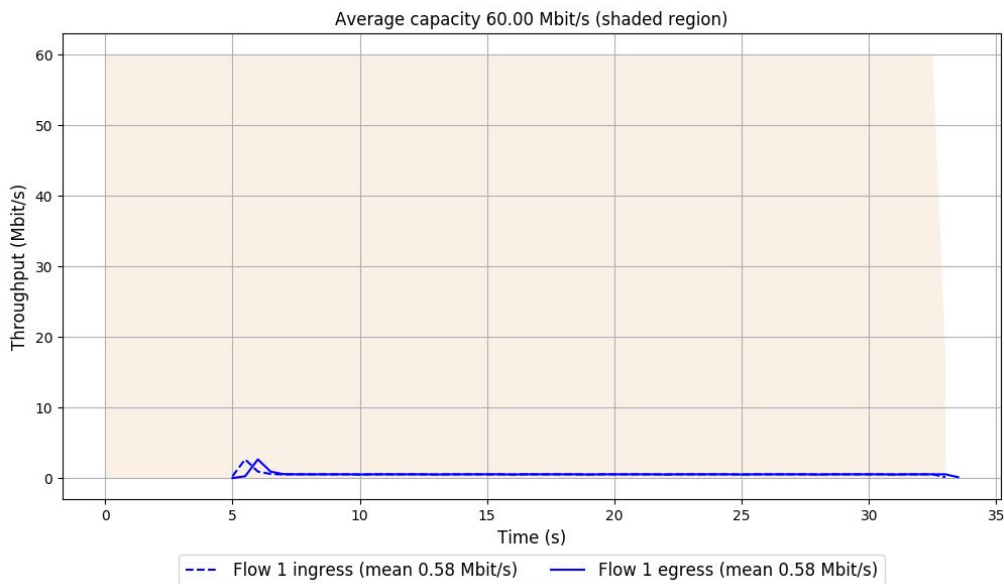


scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	0.58	0.60	0.63	3.47	3.34	3.85	0.00	0.00	0.00

Test 4: PPO [64, 32, 64]



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	0.63	129.97	0.00



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	0.58	1.35	0.00

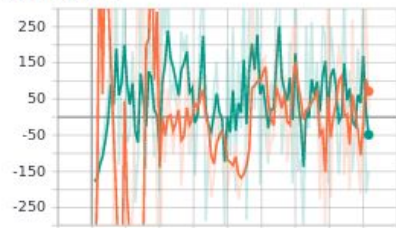
Aurora: TD3

- RL Algorithm: Twin Delayed Deep Deterministic Policy Gradient
- Architecture (of NN): [64, 32]
- History Length: 10
- Features:
- Gamma: 0.99
- Twin Delayed Deep Deterministic Policy Gradient: TD3 addresses function approximation error in Actor-Critic methods. TD3 is a direct successor of DDPG and improves it using three major tricks: clipped Q-Learning, delayed policy update, and target policy smoothing. These tricks result in substantially improved performance over baseline DDPG.

TD3 [64, 32]: Training Statistics

episode_reward

episode_reward



Orange Graph - First Checkpoint, Green Graph - Last Checkpoint

loss

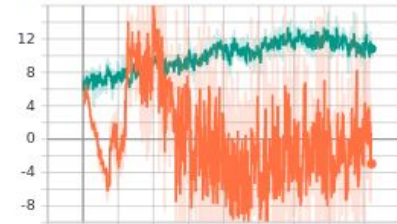
learning_rate

tag: loss/learning_rate



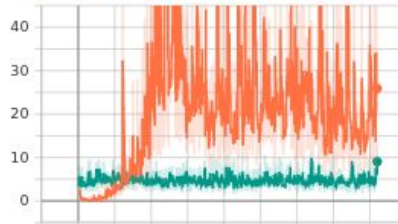
policy_loss

tag: loss/policy_loss



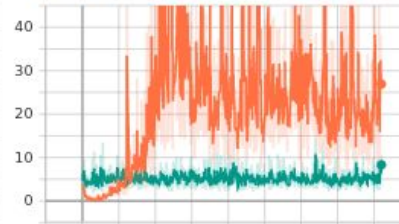
qf1_loss

tag: loss/qf1_loss



qf2_loss

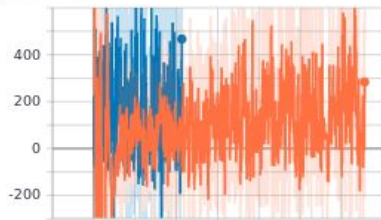
tag: loss/qf2_loss



TD3: Divergence

episode_reward

episode_reward

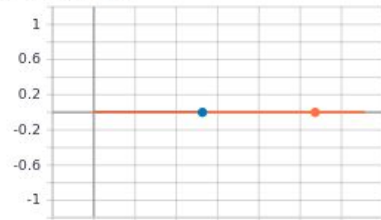


Orange Graph - First Checkpoint, Blue Graph - Last Checkpoint

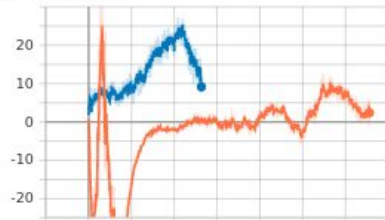


loss

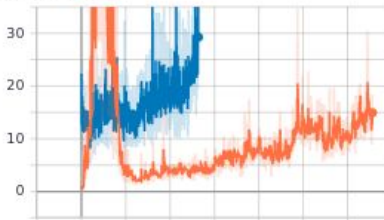
learning_rate
tag: loss/learning_rate



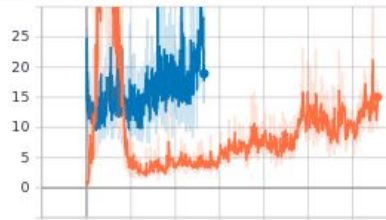
policy_loss
tag: loss/policy_loss



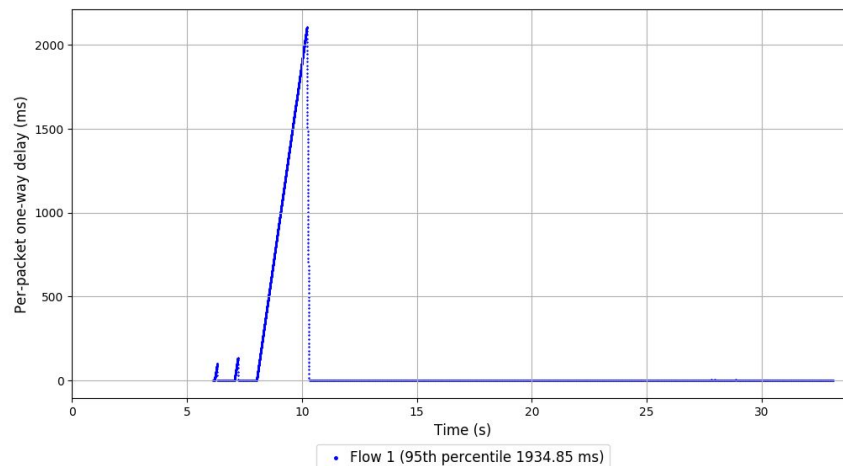
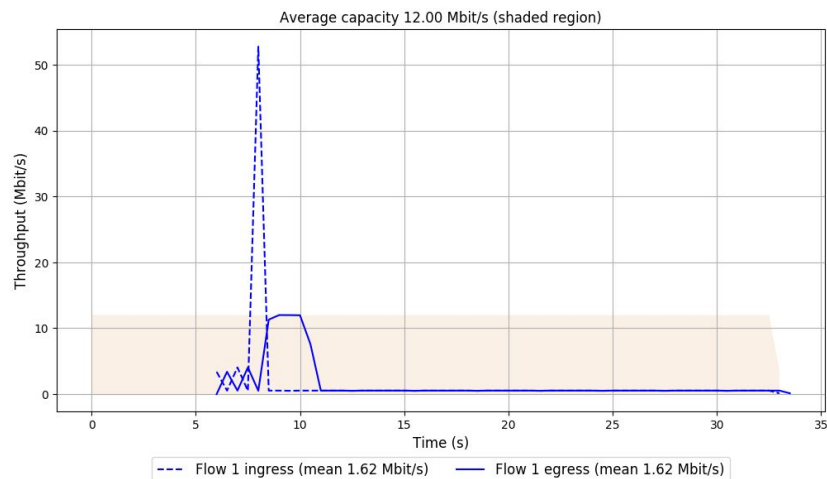
qf1_loss
tag: loss/qf1_loss



qf2_loss
tag: loss/qf2_loss

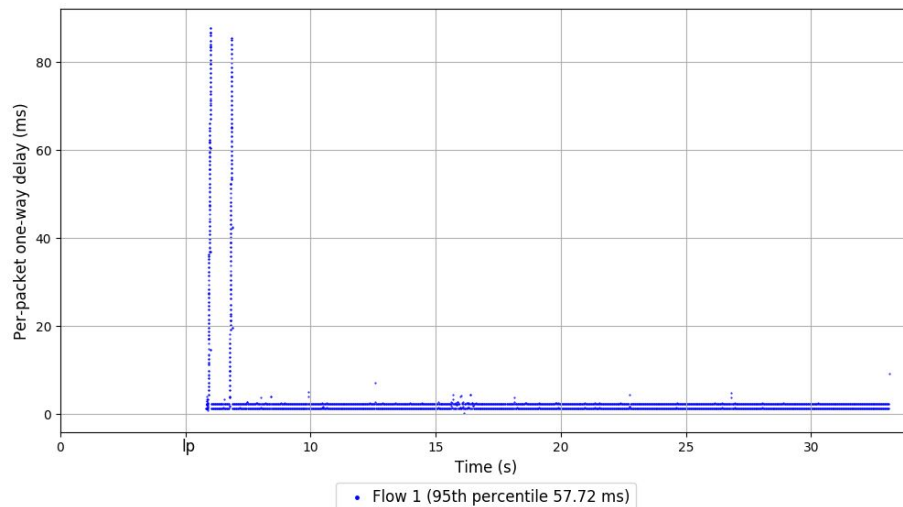
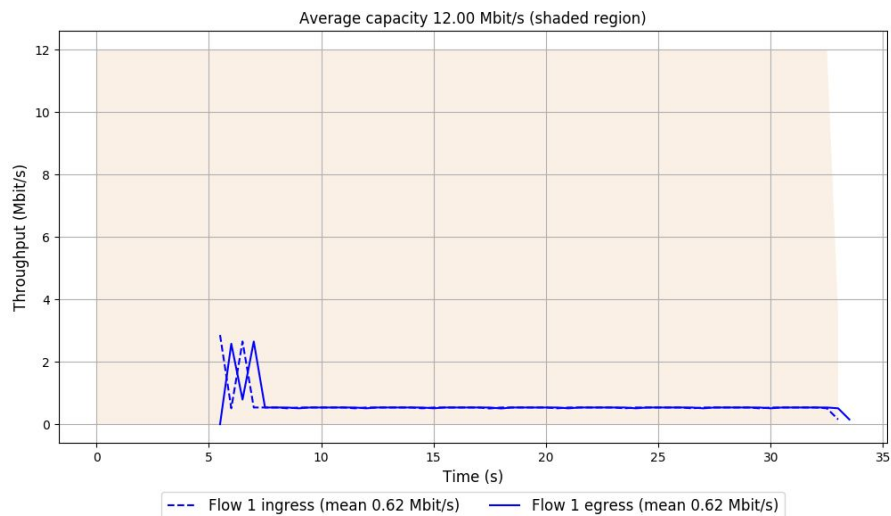


Test 1: TD3 [64, 32]



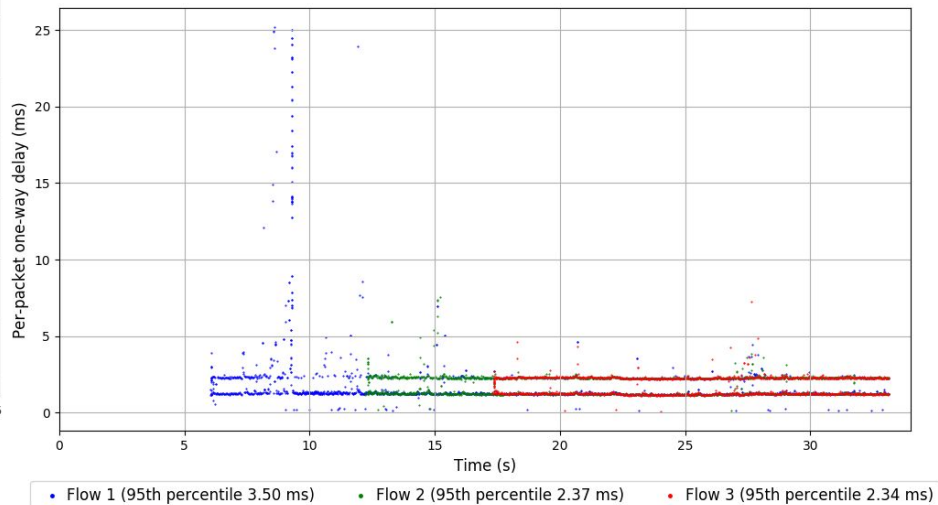
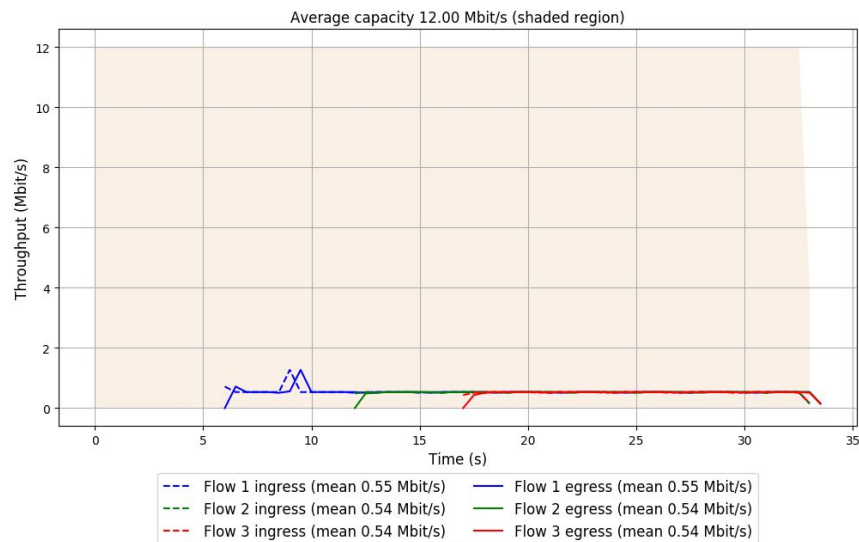
scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	1	1.62	1934.85	0.00

Test 2: TD3 [64, 32]



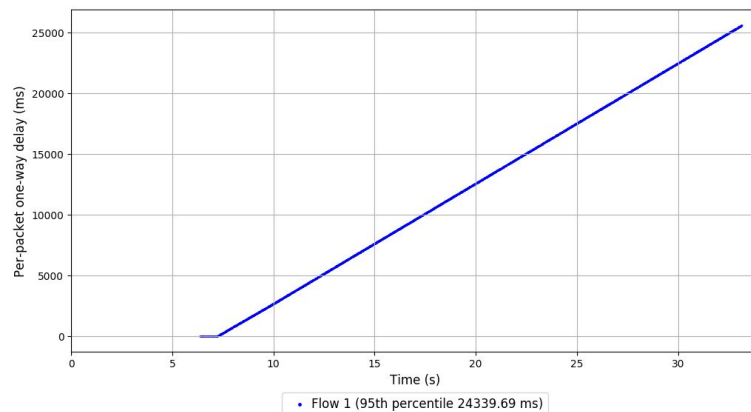
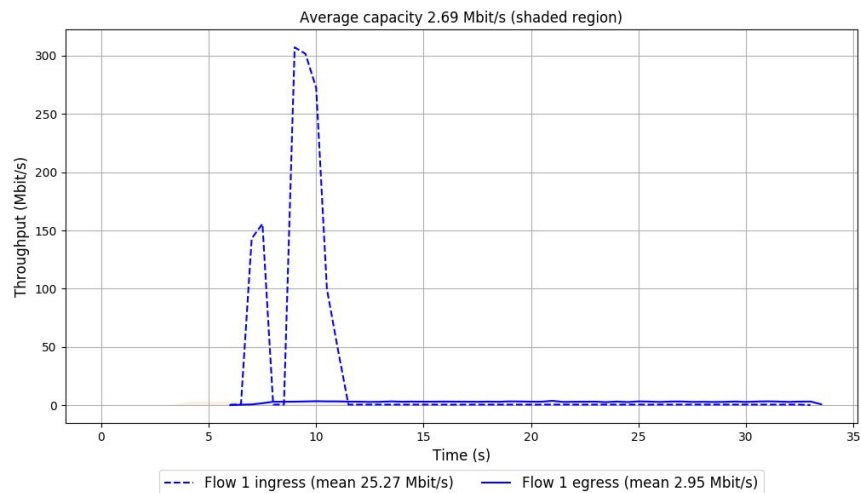
scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	5	1.81	349.16	0.27

Test 3: TD3 [64, 32]



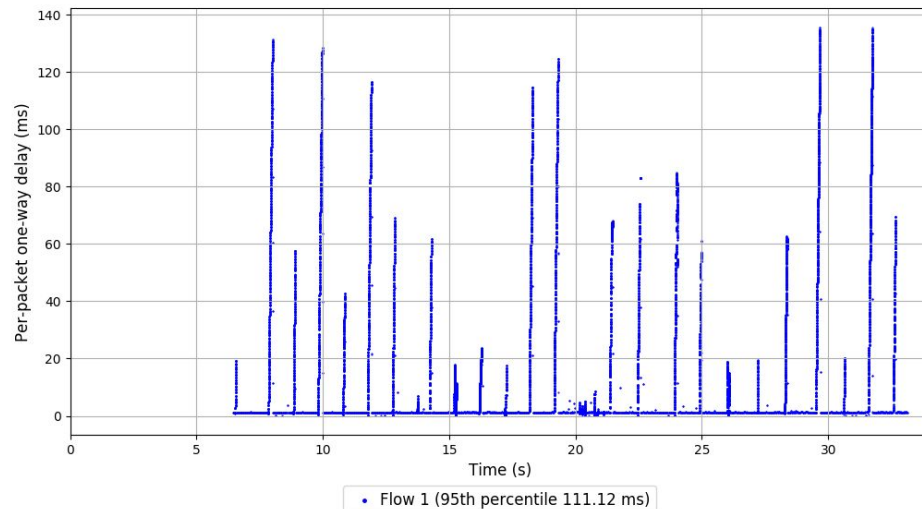
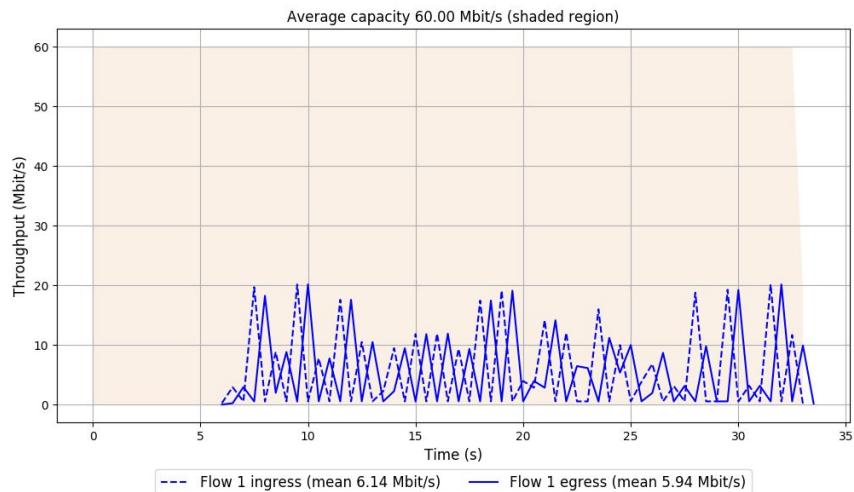
scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	0.55	0.54	0.54	3.50	2.37	2.34	0.00	0.00	0.15

Test 4: TD3 [64, 32]



scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	2.53	1.02	9.76	2820.67	7738.59	8570.39	6.41	95.32	32.84

Test 5: TD3 [64,32]



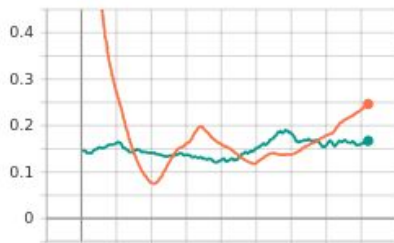
scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	1	5.94	111.12	3.28

Aurora: SAC

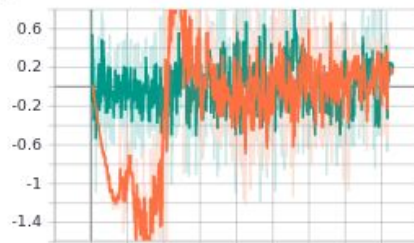
- RL Algorithm: Soft Actor Critic
- Architecture (of NN): [64, 32]
- History Length: 10
- Features:
- Gamma: 0.99
- Soft Actor Critic: SAC is an off-policy maximum entropy deep reinforcement learning with a stochastic actor. SAC incorporates the double Q-learning trick from TD3. A key feature of SAC, and a major difference with common RL algorithms, is that it is trained to maximise a trade-off between expected return and entropy, a measure of randomness in the policy.

SAC [64, 32]: Training Statistics

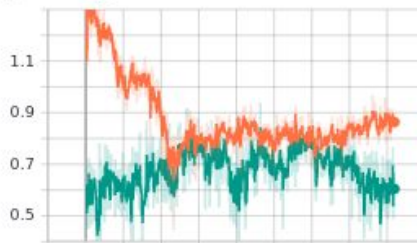
ent_coef
tag: loss/ent_coef



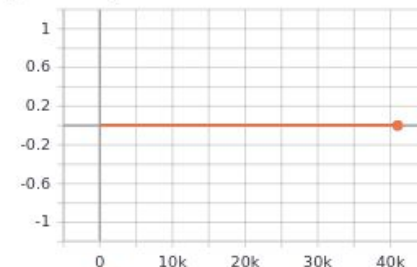
ent_coef_loss
tag: loss/ent_coef_loss



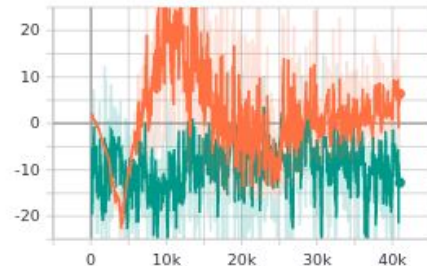
entropy
tag: loss/entropy



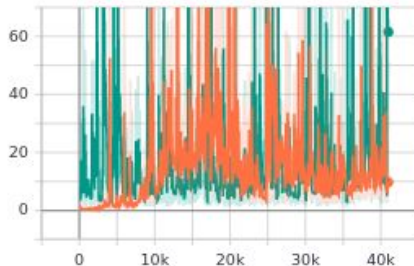
learning_rate
tag: loss/learning_rate



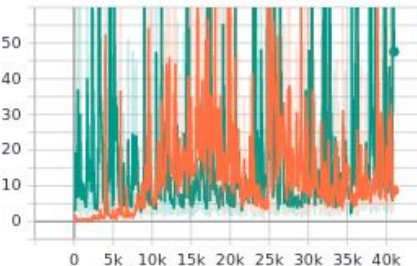
policy_loss
tag: loss/policy_loss



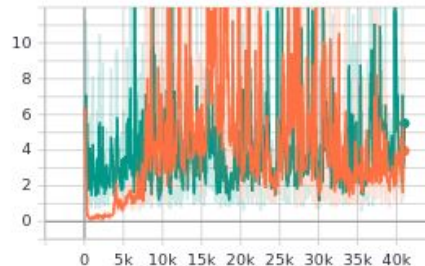
qf1_loss
tag: loss/qf1_loss



qf2_loss
tag: loss/qf2_loss

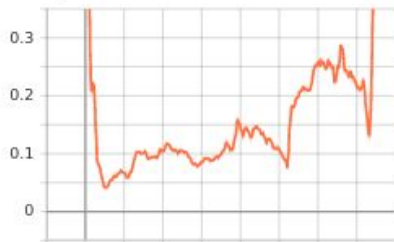


value_loss
tag: loss/value_loss

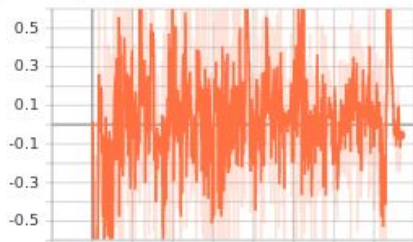


SAC: Divergence

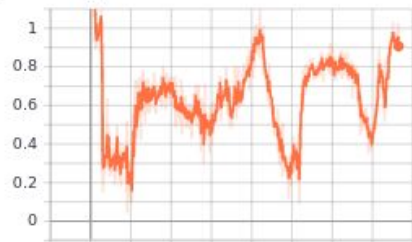
ent_coef
tag: loss/ent_coef



ent_coef_loss
tag: loss/ent_coef_loss



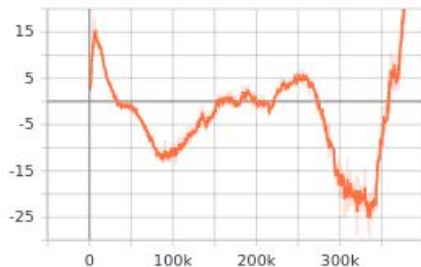
entropy
tag: loss/entropy



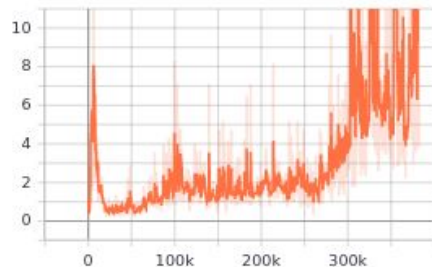
learning_rate
tag: loss/learning_rate



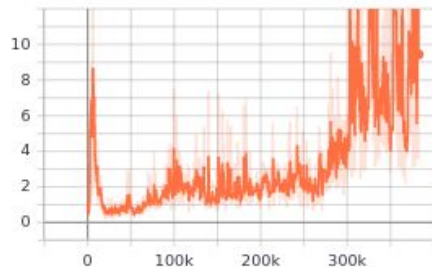
policy_loss
tag: loss/policy_loss



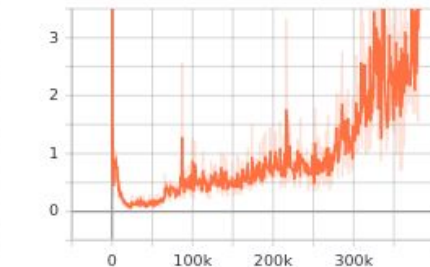
qf1_loss
tag: loss/qf1_loss



qf2_loss
tag: loss/qf2_loss

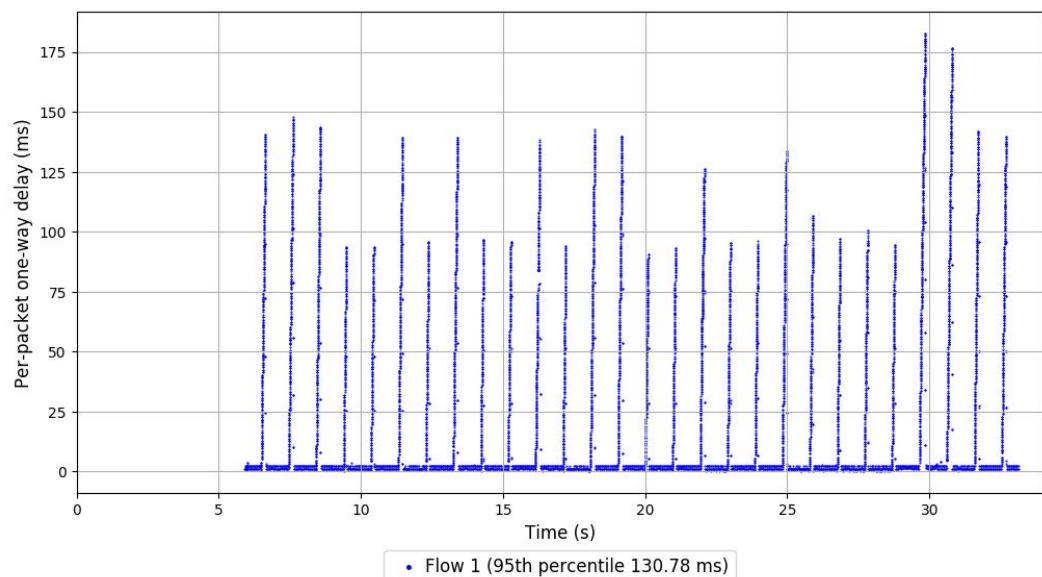
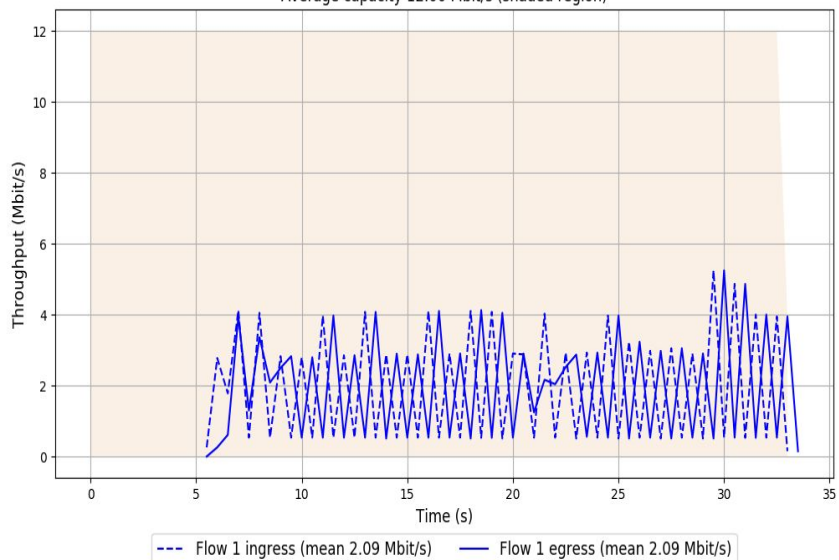


value_loss
tag: loss/value_loss



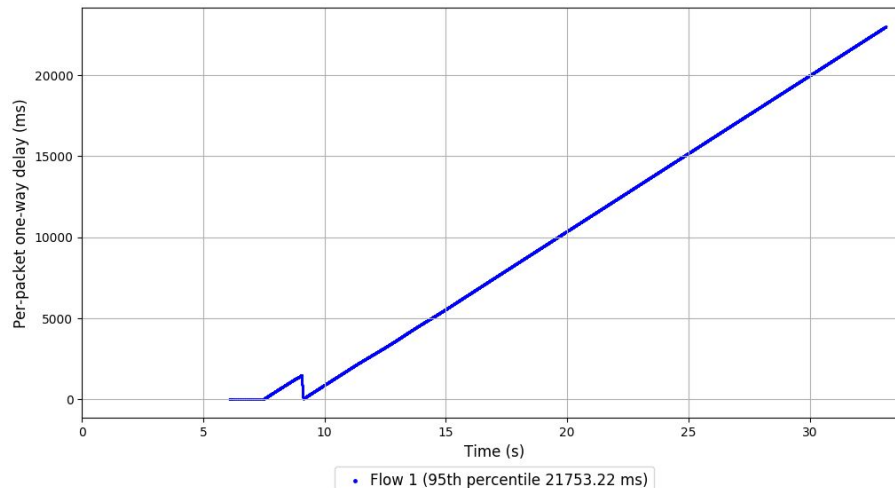
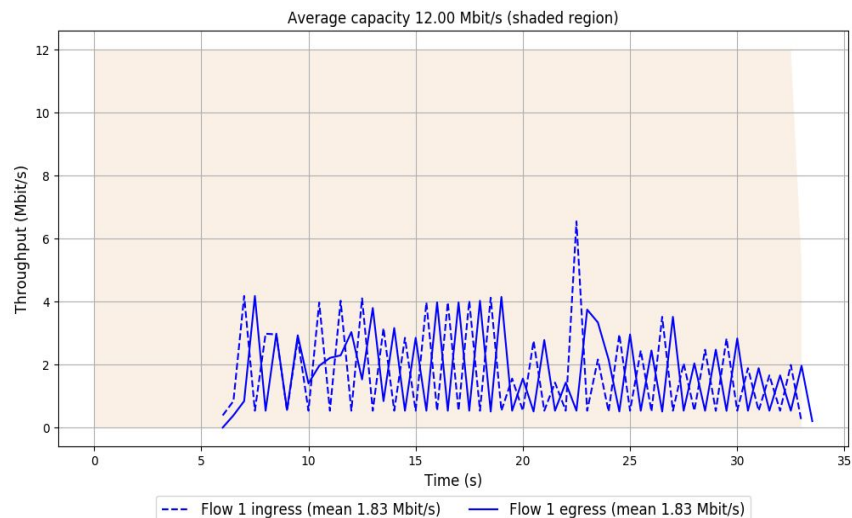
Test 1: SAC [64, 32]

Average capacity 12.00 Mbit/s (shaded region)



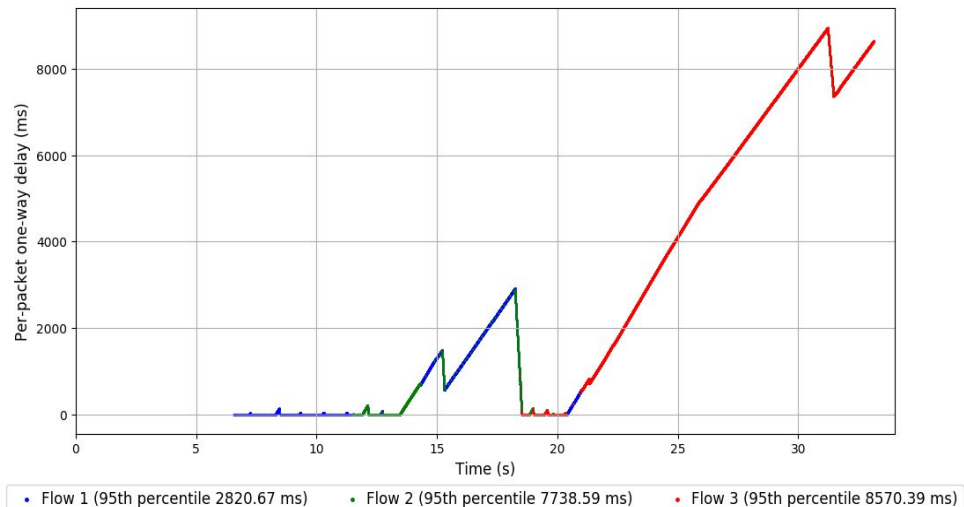
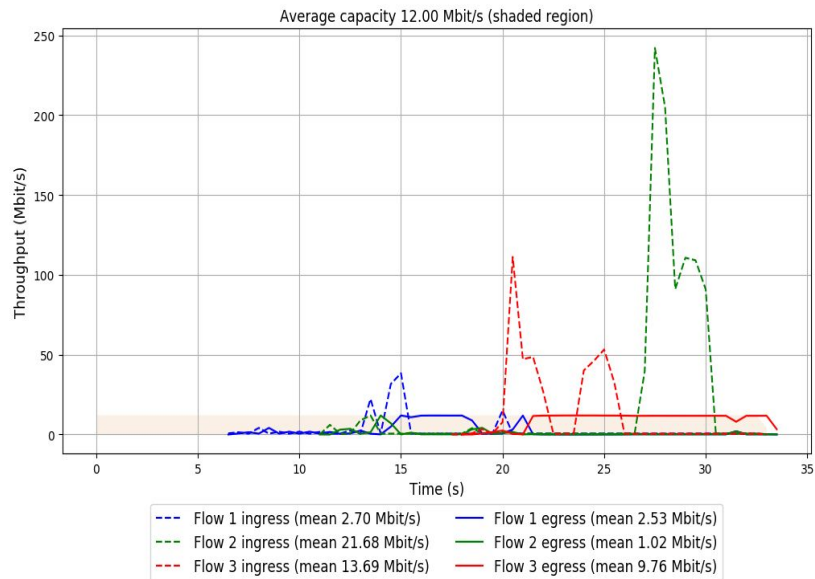
scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	1	2.09	130.78	0.00

Test 2: SAC [64, 32]



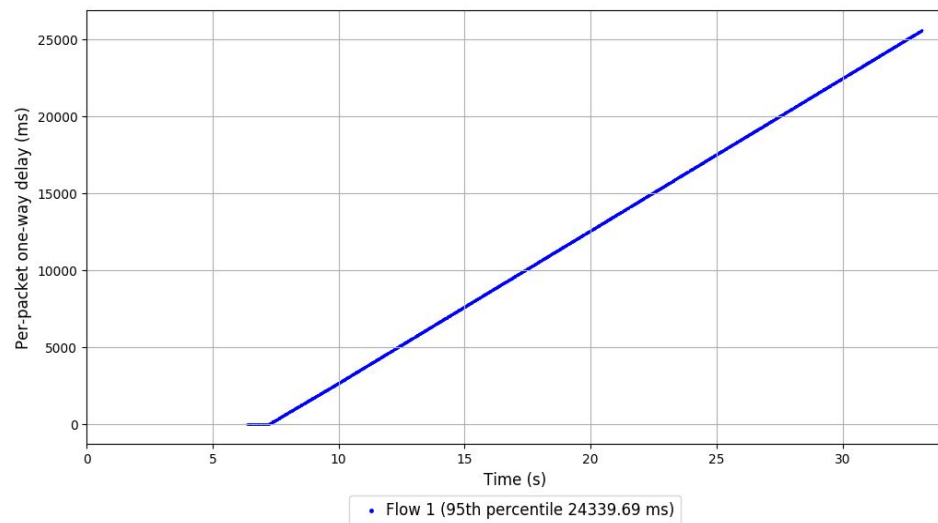
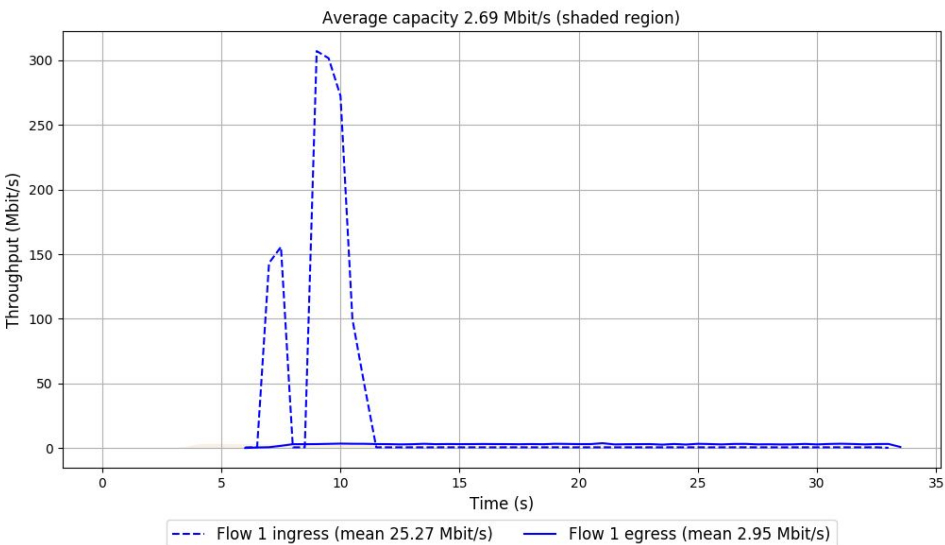
scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	5	3.96	4500.13	16.24

Test 3: SAC [64, 32]



scheme	# runs	mean avg tput (Mbit/s)			mean 95th-%ile delay (ms)			mean loss rate (%)		
		flow 1	flow 2	flow 3	flow 1	flow 2	flow 3	flow 1	flow 2	flow 3
PCC-RL	1	2.53	1.02	9.76	2820.67	7738.59	8570.39	6.41	95.32	32.84

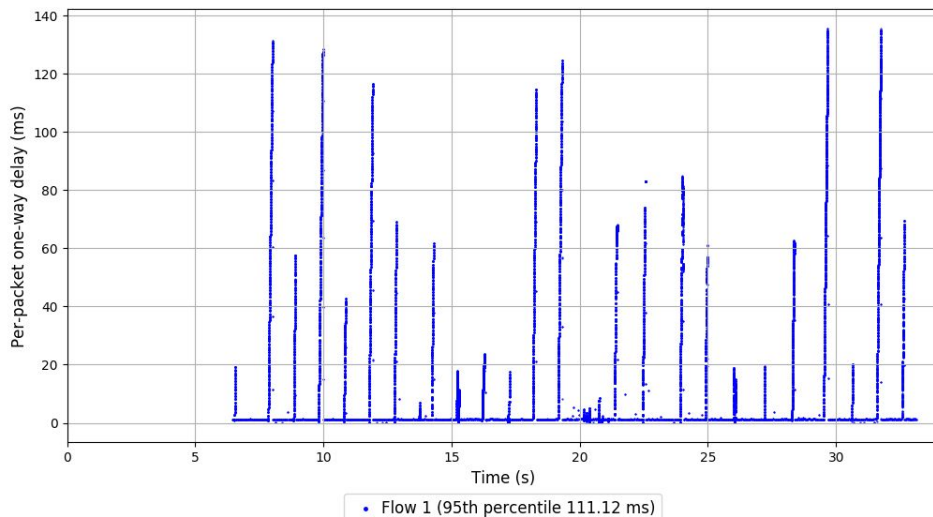
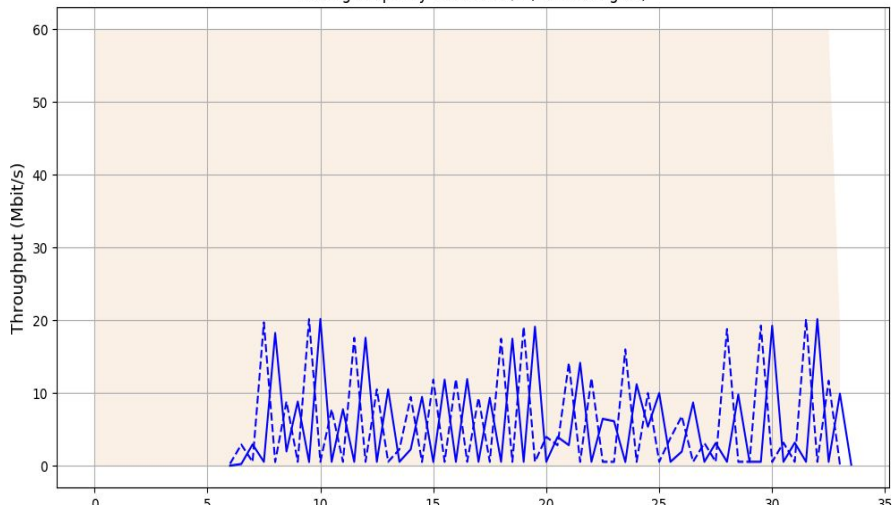
Test 4: SAC [64, 32]



scheme	# runs	mean avg tput (Mbit/s)	mean 95th-%ile delay (ms)	mean loss rate (%)
		flow 1	flow 1	flow 1
PCC-RL	1	2.95	24339.69	88.31

Test 5: SAC [64,32]

Average capacity 60.00 Mbit/s (shaded region)



scheme	# runs	mean avg tput (Mbit/s) flow 1	mean 95th-%ile delay (ms) flow 1	mean loss rate (%) flow 1
PCC-RL	1	5.94	111.12	3.28

Multi-Objective: Reward Engineering

Reward is a linear combination of throughput, latency, loss

How does the performance change with the individual weights

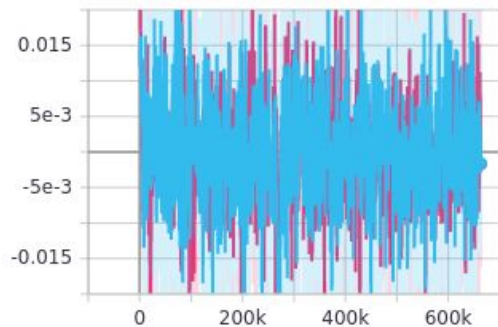
Reward in earlier experiments: $10 \cdot \text{throughput} - 1000 \cdot \text{latency} - 2000 \cdot \text{loss}$

High throughput (expected): $20 \cdot \text{throughput} - 1000 \cdot \text{latency} - 2000 \cdot \text{loss}$

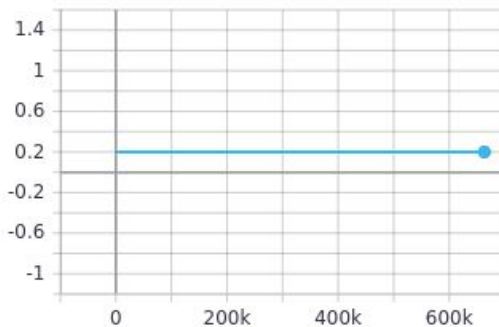
Low latency (expected): $5 \cdot \text{throughput} - 1000 \cdot \text{latency} - 2000 \cdot \text{loss}$

PPO [64, 32] HighT: Training Statistics

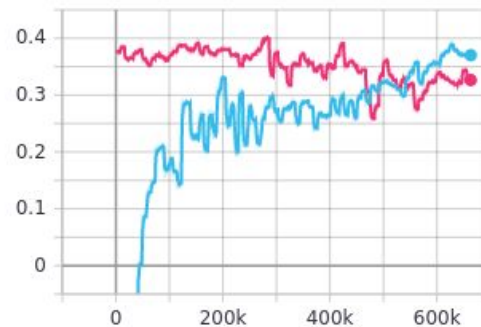
advantage
tag: input_info/advantage



clip_range
tag: input_info/clip_range



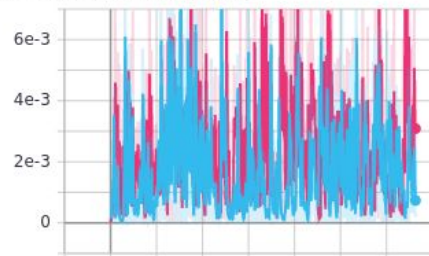
discounted_rewards
tag: input_info/discounted_rewards



Blue graph - First Checkpoint, Pink graph - Last Checkpoint

PPO [64, 32] HighT: Training Statistics

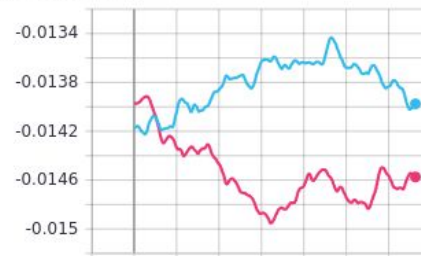
approximate_kullback-leibler
tag: loss/approximate_kullback-leibler



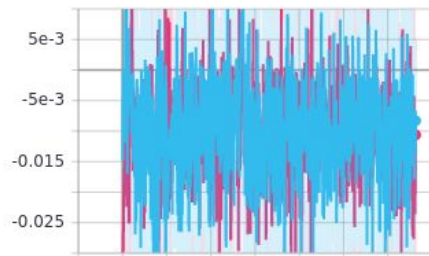
clip_factor
tag: loss/clip_factor



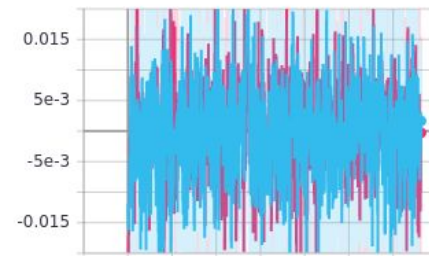
entropy_loss
tag: loss/entropy_loss



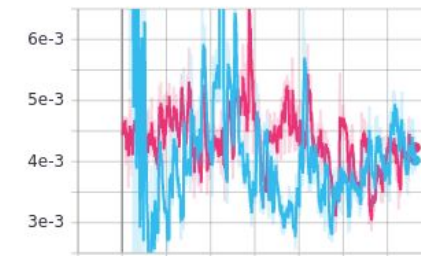
loss
tag: loss/loss



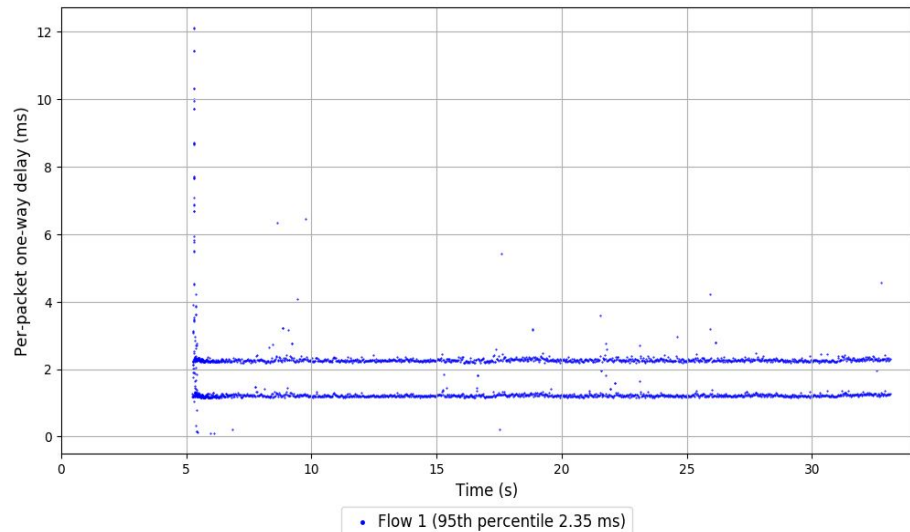
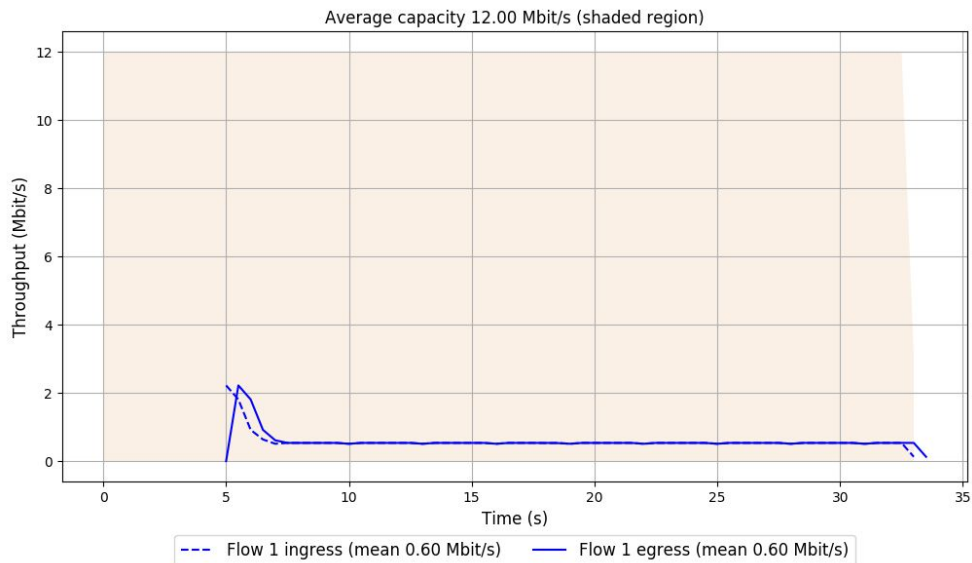
policy_gradient_loss
tag: loss/policy_gradient_loss



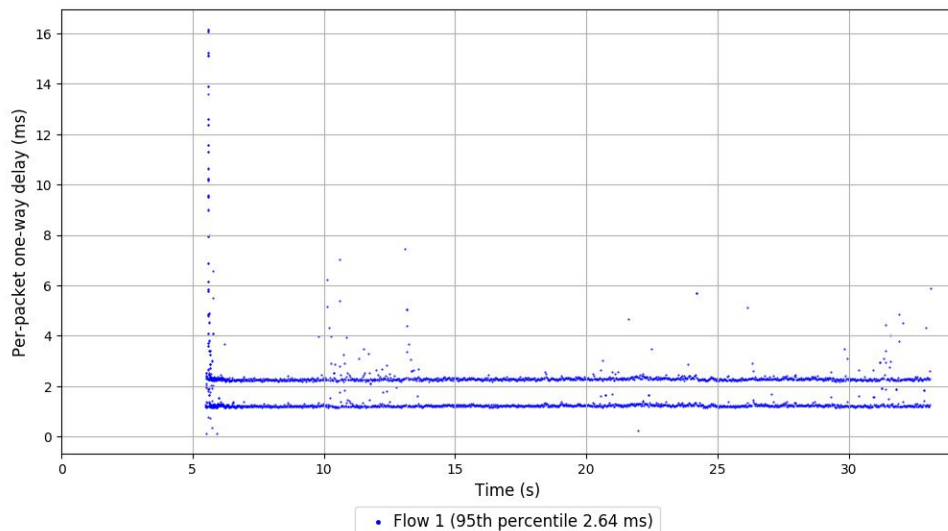
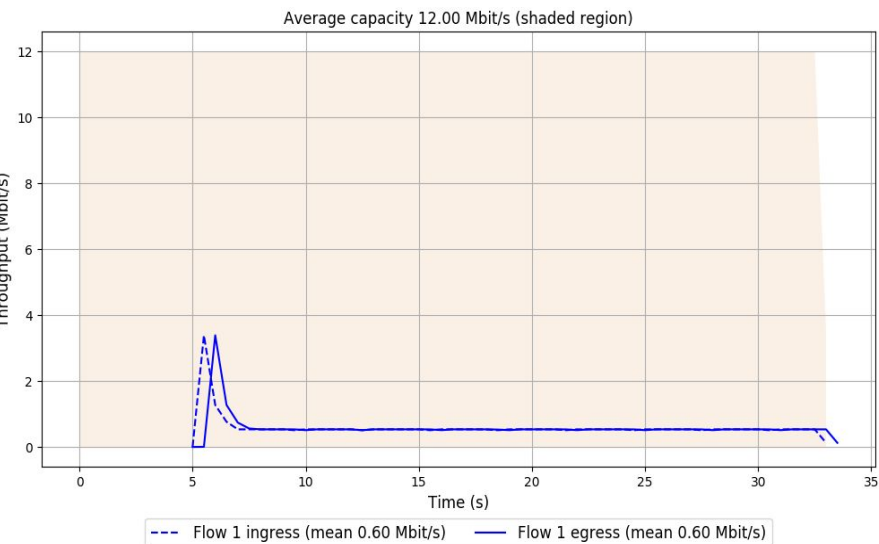
value_function_loss
tag: loss/value_function_loss



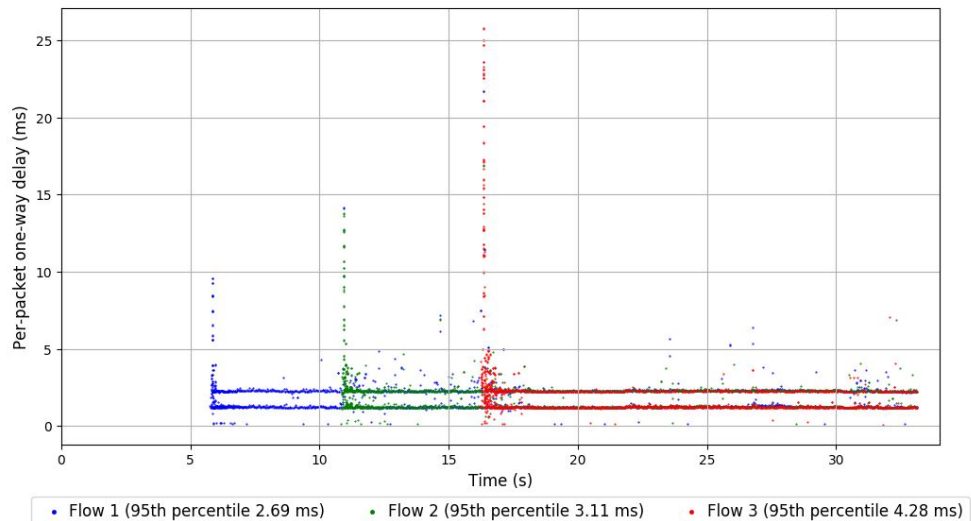
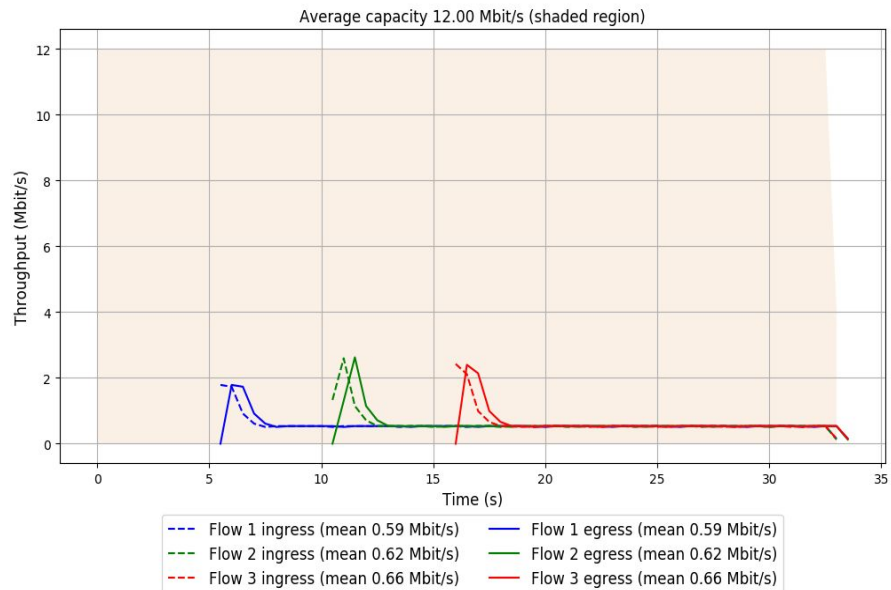
Test 1: PPO [64, 32] HighT



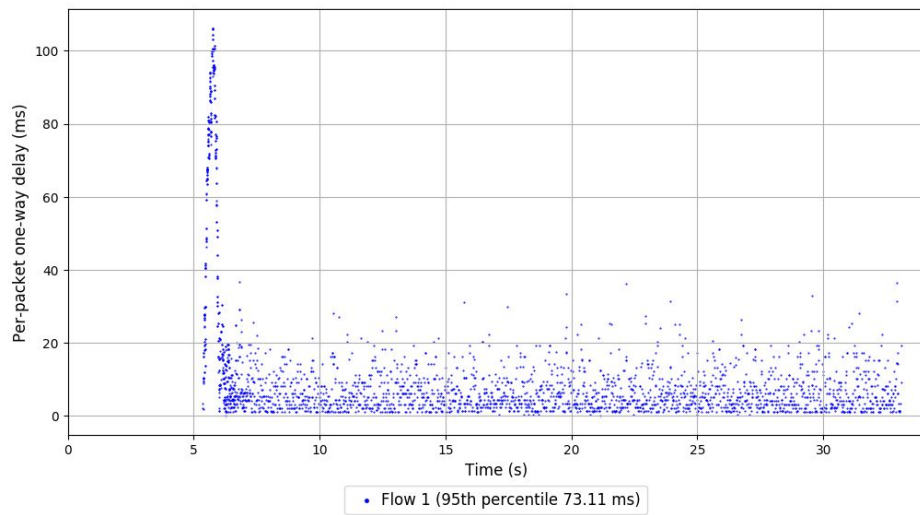
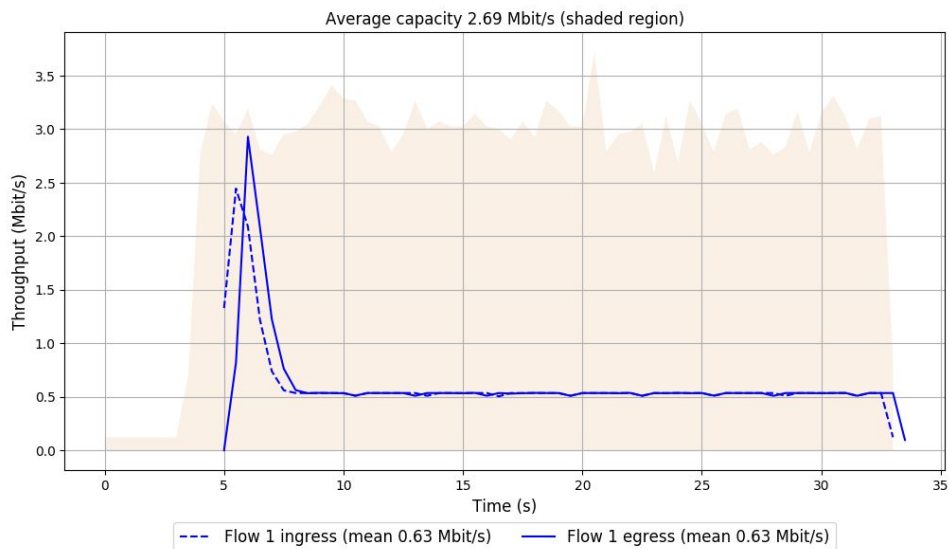
Test 2: PPO [64, 32] HighT



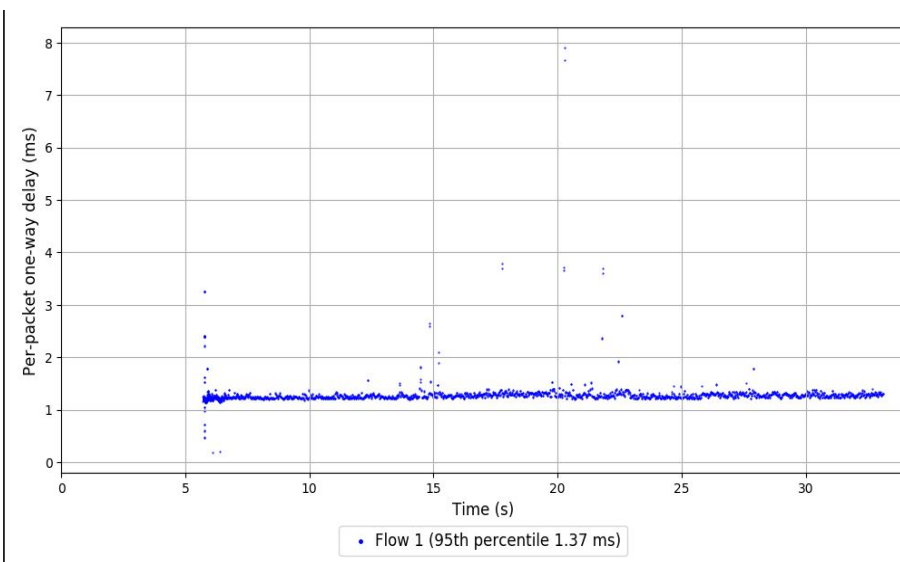
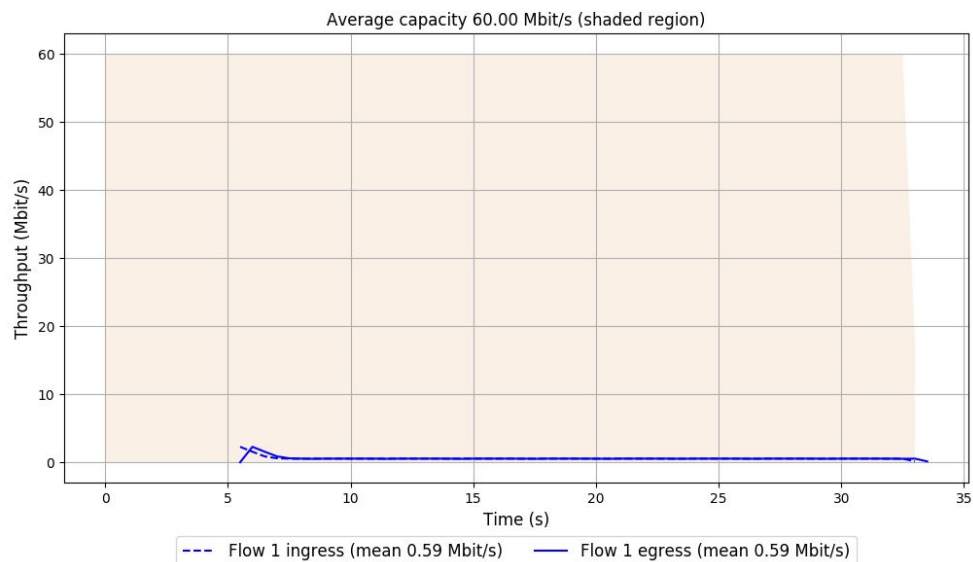
Test 3: PPO [64, 32] HighT



Test 4: PPO [64, 32] HighT

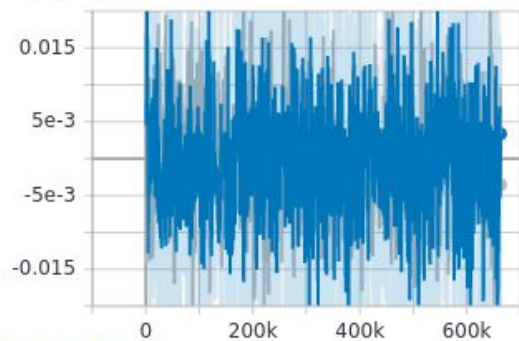


Test 5: PPO [64,32] HighT

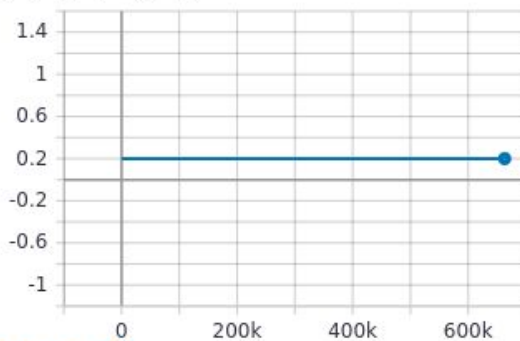


PPO [64, 32] LowLat: Training Statistics

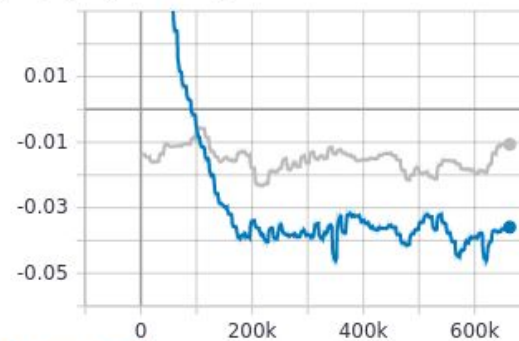
advantage
tag: input_info/advantage



clip_range
tag: input_info/clip_range



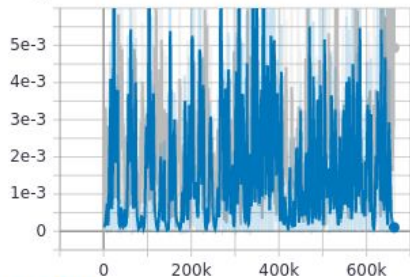
discounted_rewards
tag: input_info/discounted_rewards



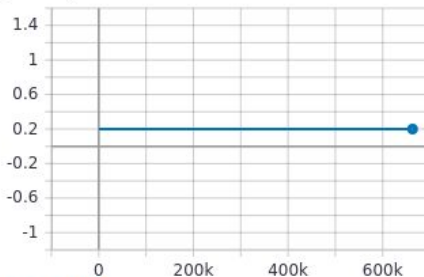
Blue Graph - First Checkpoint, Grey Graph - Last Checkpoint

PPO [64, 32] LowLat: Training Statistics

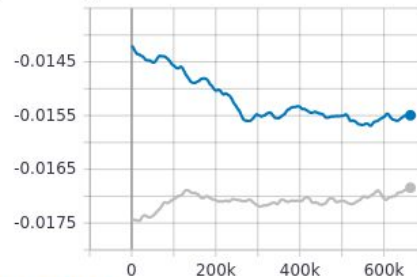
approximate_kullback-leibler
tag: loss/approximate_kullback-leibler



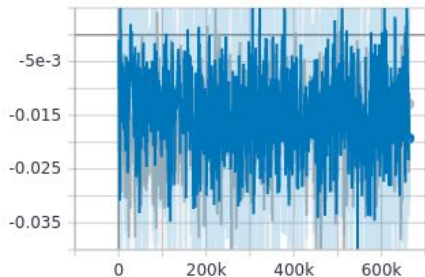
clip_factor
tag: loss/clip_factor



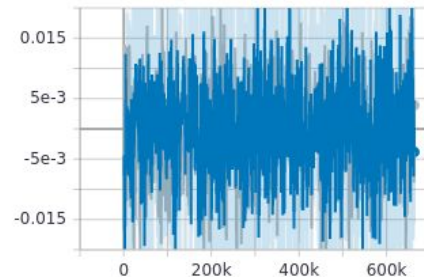
entropy_loss
tag: loss/entropy_loss



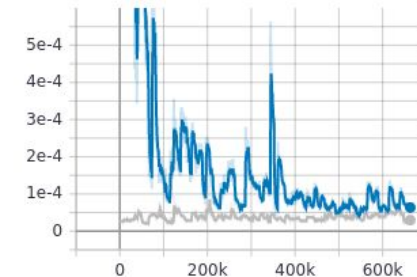
loss
tag: loss/loss



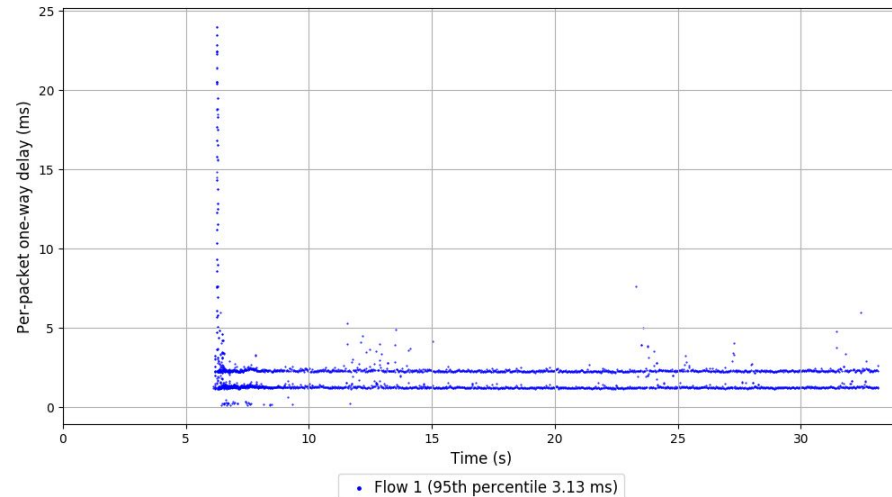
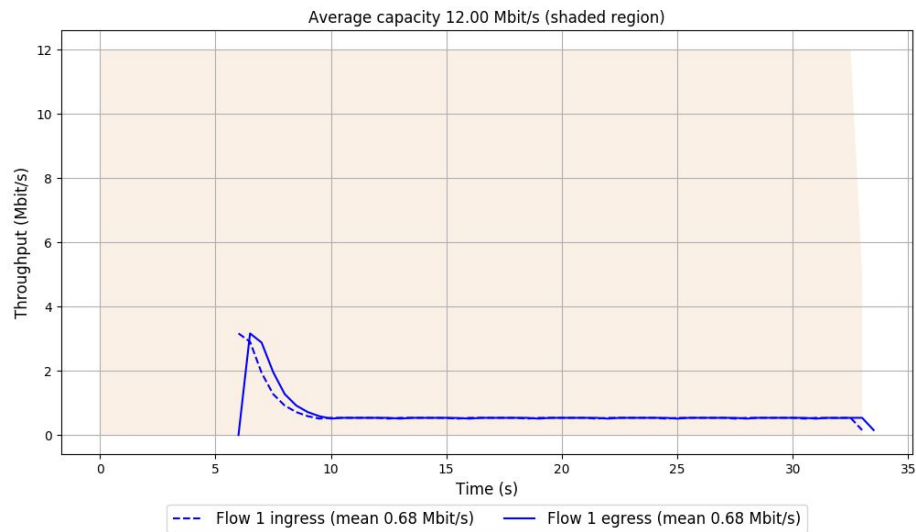
policy_gradient_loss
tag: loss/policy_gradient_loss



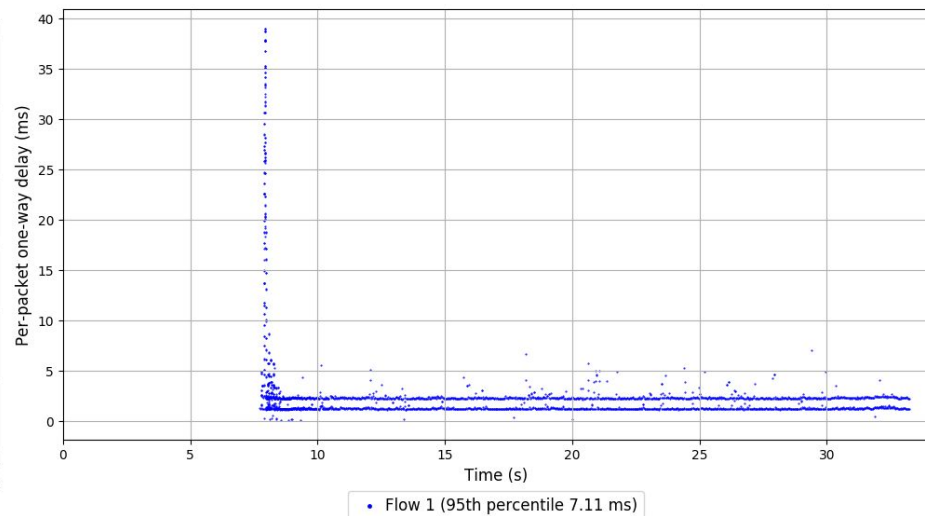
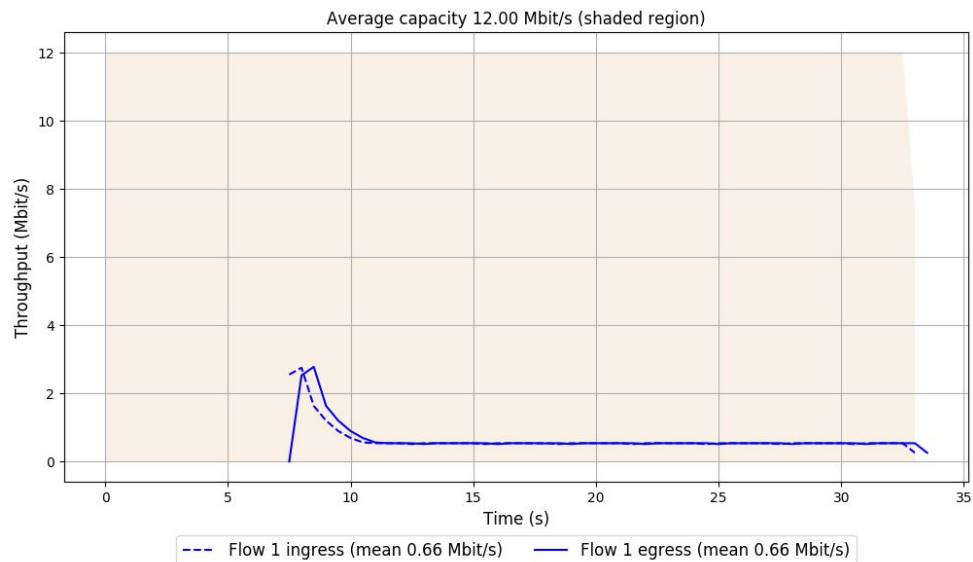
value_function_loss
tag: loss/value_function_loss



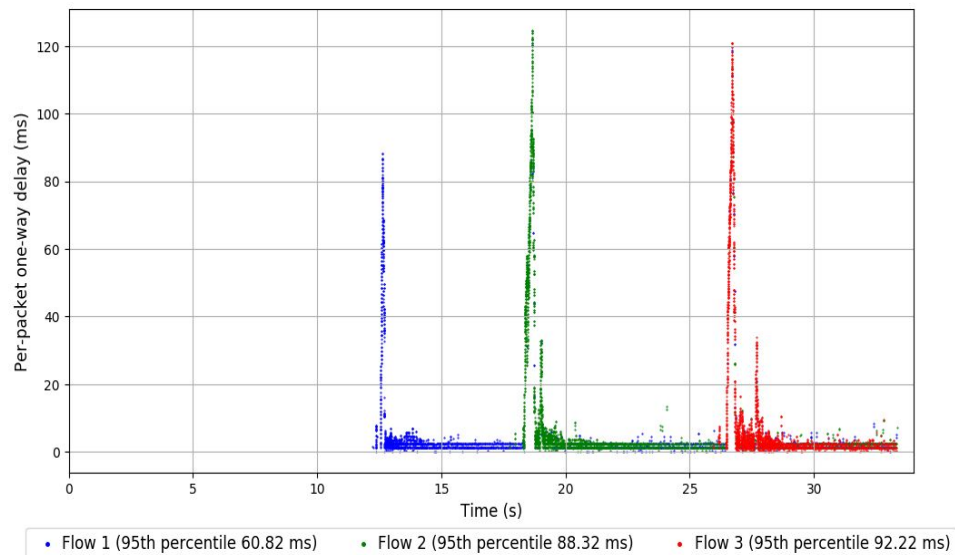
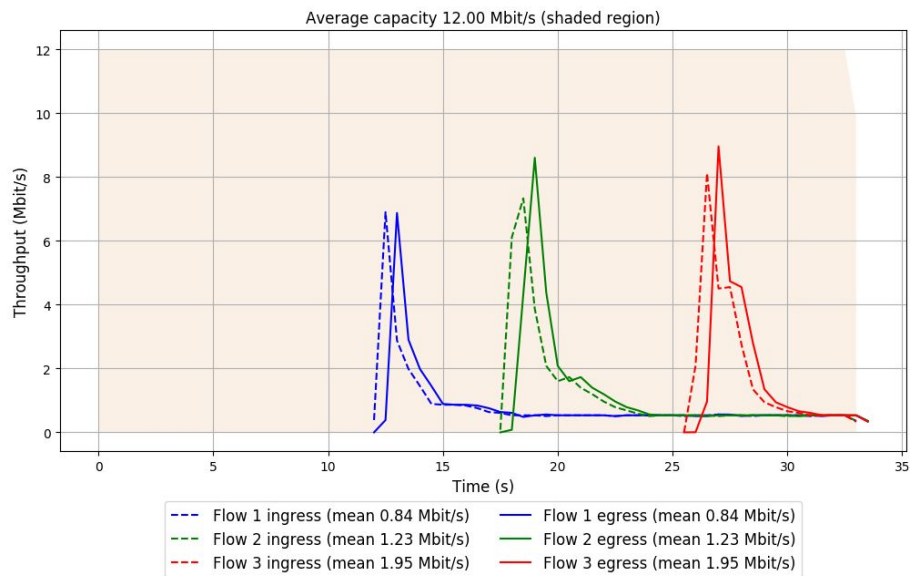
Test 1: PPO [64, 32] LowLat



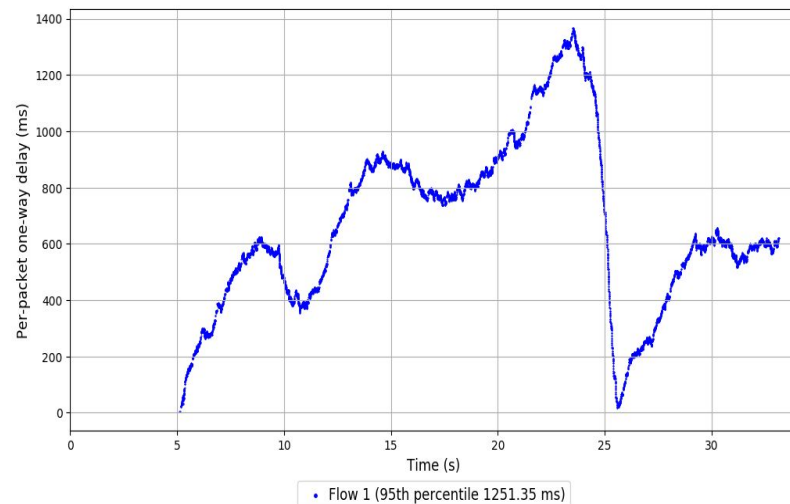
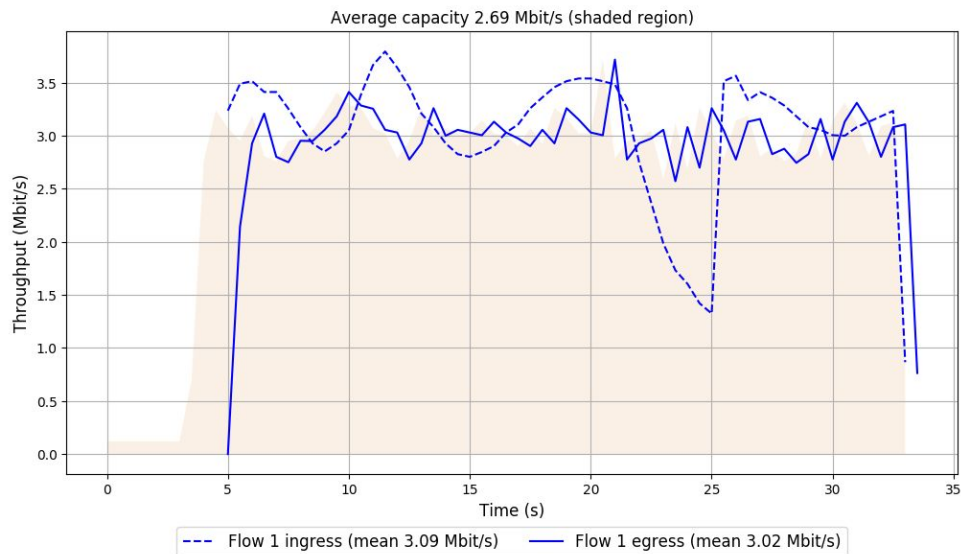
Test 2: PPO [64, 32] LowLat



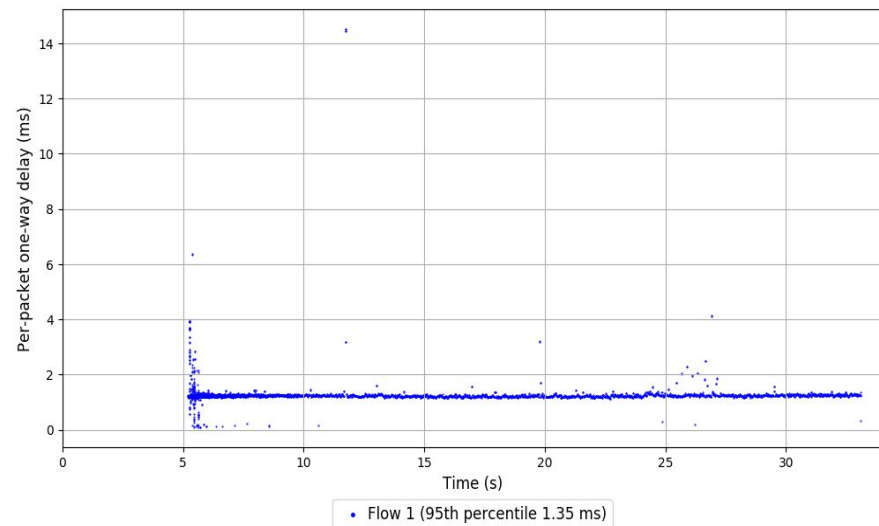
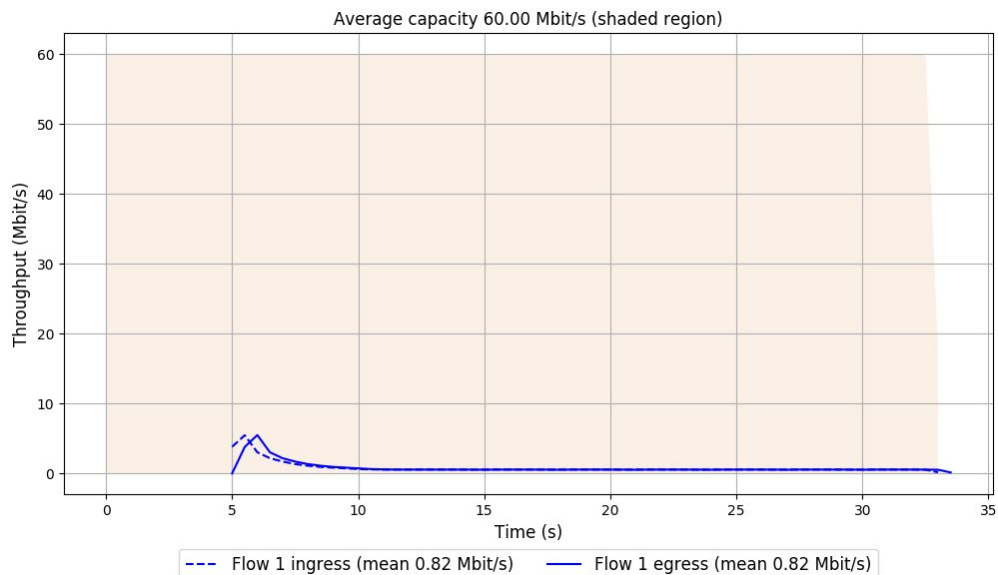
Test 3: PPO [64, 32] LowLat



Test 4: PPO [64, 32] LowLat



Test 5: PPO [64,32] LowLat



Roadblocks

- Code/libraries used are outdated
 - Python2 -> Python3
 - Tensorflow 1.14 -> Tensorflow 2+
- Outdated kernel issues with Orca and DeepCC
- All models trained on CPU

Further Avenues

- Model-based approaches
- Meta-Learning
- Reward Engineering
- Competitive Learning

References

- [1] Nathan Jay, Noga Rotman, Brighten Godfrey, Michael Schapira, Aviv Tamar. A Deep Reinforcement Learning Perspective on Internet Congestion Control. *Proceedings of the 36th International Conference on Machine Learning*, PMLR 97:3050-3059, 2019.
- [2] Yiqing Ma, Han Tian, Xudong Liao, Junxue Zhang, Weiyan Wang, Kai Chen, and Xin Jin. 2022. Multi-objective congestion control. In *Proceedings of the Seventeenth European Conference on Computer Systems (EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 218–235. <https://doi.org/10.1145/3492321.3519593>
- [3] Francis Y. Yan, Jestin Ma, Greg D. Hill, Deepti Raghavan, Riad S. Wahby, Philip Levis, and Keith Winstein. 2018. Pantheon: the training ground for internet congestion-control research. In *Proceedings of the 2018 USENIX Conference on Usenix Annual Technical Conference (USENIX ATC '18)*. USENIX Association, USA, 731–743.
- [4] <https://stable-baselines.readthedocs.io/en/master/index.html>
- [5] <https://spinningup.openai.com/en/latest/index.html>