The Problem

The DVRP aims to find the least cost routes to serve a set of stochastic customers with a given number of capacitated vehicles

Objective

Evaluate performance of RL algorithms against heursitic and metaheuristic algorithms - measured by solution quality, computation time and dynamic adaption

Dynamic Vehicle Routing Problem

Reinforcement Learning vs
Traditional Methods

Methodology

- Evaluated on adapted CVRP Solomon
 Instances with random, clustered and random-clustered instances ranging from 25-100 customers
- Benchmark performance obtained from PYVRPs Genetic Algorithm Solver

Algorithms

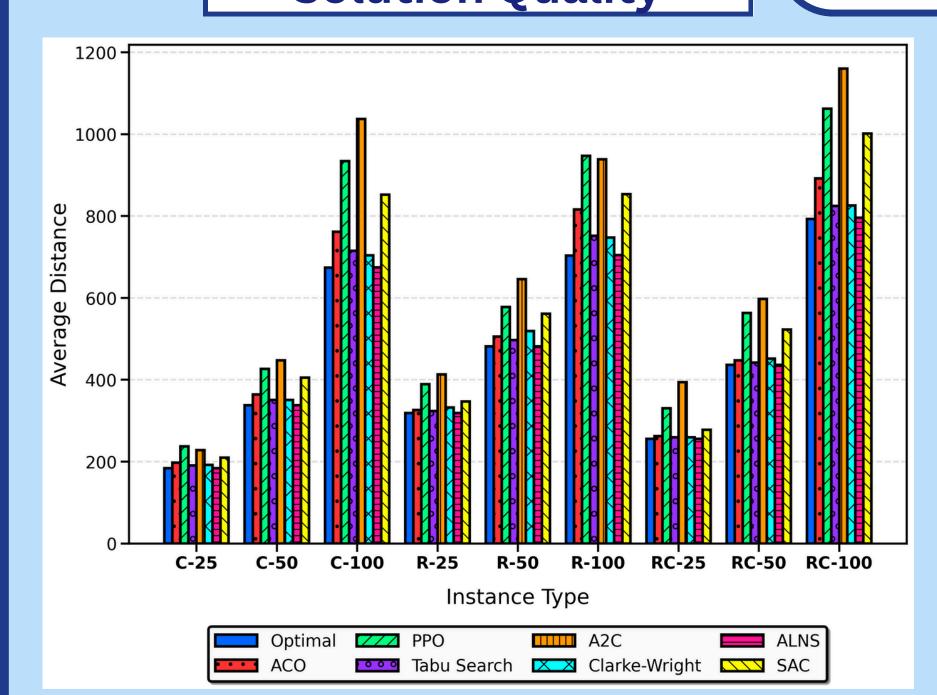
Traditional:

- Clarke-Wright Savings
- Tabu Search
- Ant ColonyOptimisation (ACO)
- Adaptive Large
 Neighbourhood Search
 (ALNS)

RL:

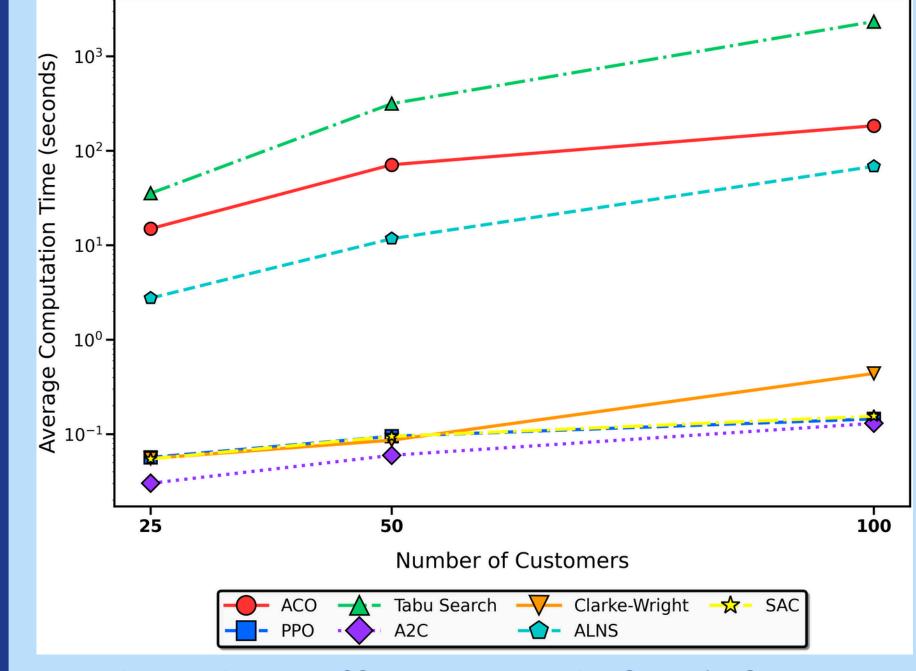
- Soft Actor-Critic (SAC)
- Advantage Actor-Critic (A2C)
- Proximal PolicyOptimisation (PPO)

Solution Quality Results



Traditional optimisation methods **consistently outperform** RL approaches. ALNS achieves near-optimal solutions (0.05% gap), followed by Clarke-Wright (4.26%) and ACO (4-20%). RL methods (PPO, A2C, SAC) show poor performance with 12-41% optimality gaps. All methods adapted to **dynamic customers**.

Computation Time



RL algorithms offer extremely **fast inference** (0.03-0.16s) but require **substantial upfront training**. Clarke-Wright provides quick solutions (0.05-0.47s). ALNS balances quality and speed (2.4-79.8s). ACO and Tabu Search deliver solutions at significantly longer runtimes (15-2400+s).

Conclusions

Current RL architectures are **fundamentally limited** for constrained vehicle routing problems, despite their ability to **adapt dynamically**. Traditional methods' superiority in this research demonstrates their continued effectiveness and highlights the need for **innovations** to render RL viable in this domain



Ryan Schapiro: schrya010@myuct.ac.za
Jayden Moore: mrxjay001@myuct.ac.za
Aidan Brand: brnaid002@myuct.ac.za

Supervised by Krupa Prag: krupa.prag@myuct.ac.za

