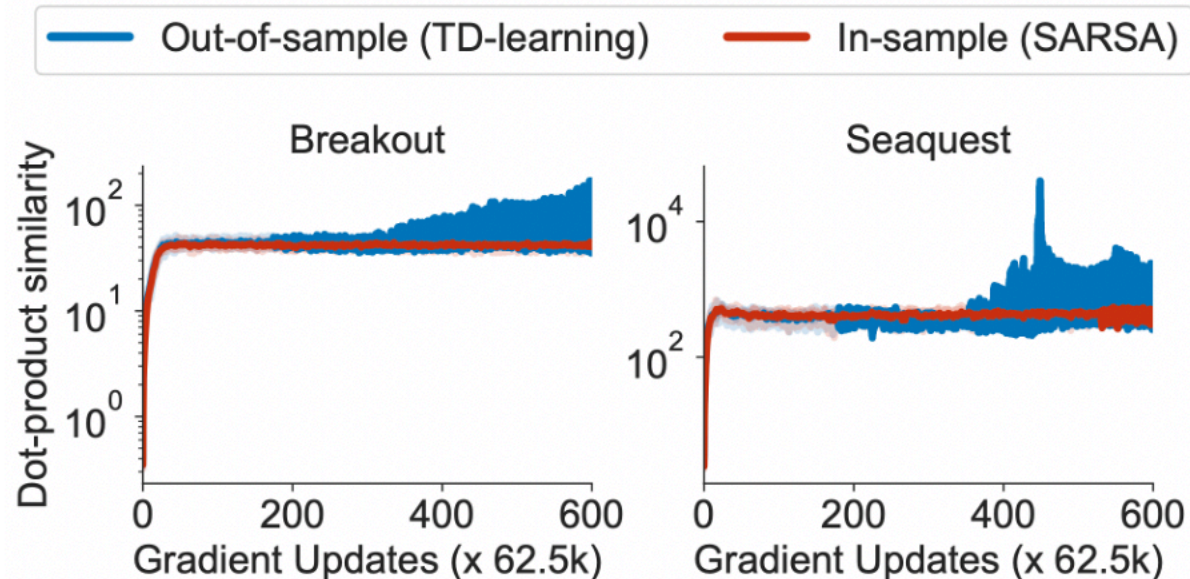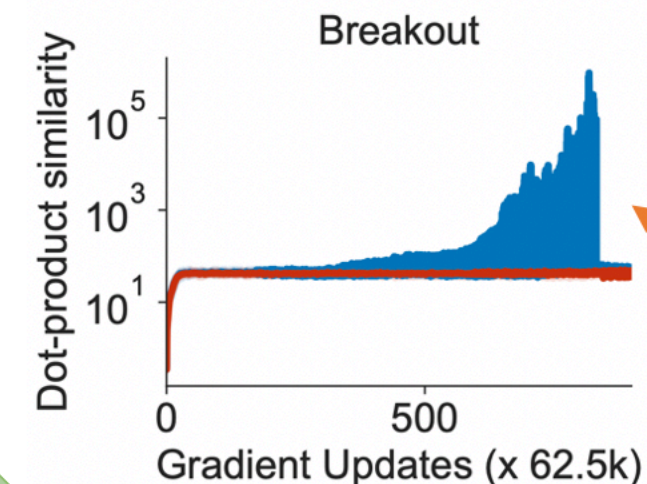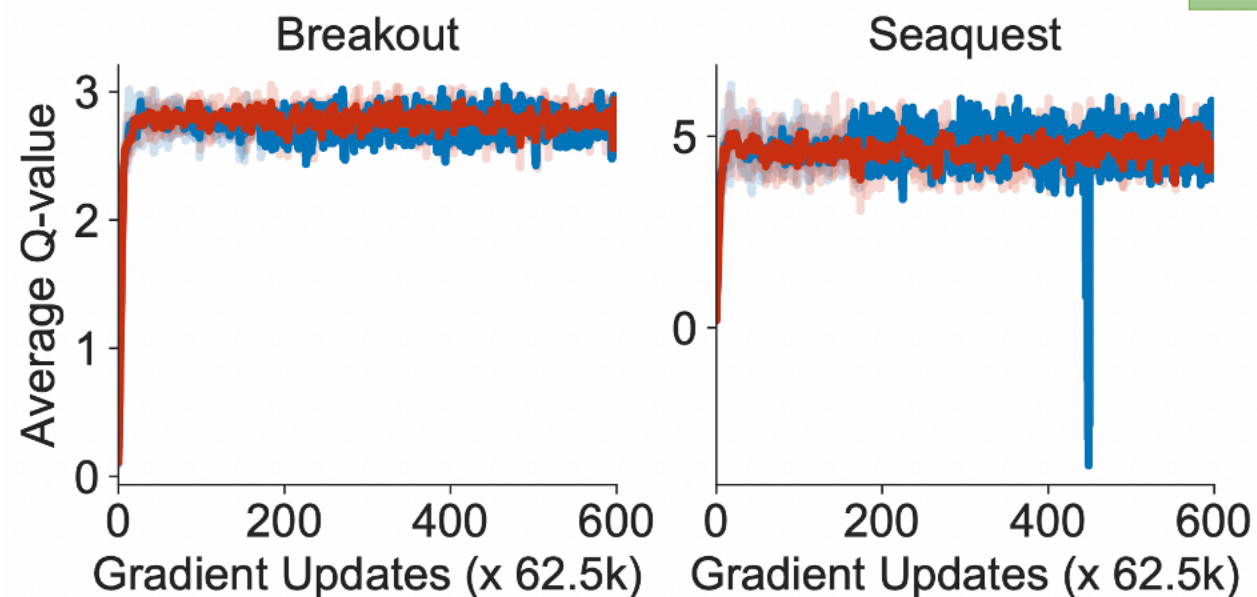Large feature dot products arise when out-of-sample actions are used in TD-learning compared to SARSA, despite similar Q-values

Large feature dot products <u>eventually</u> imply divergent Q-values

— Out-of-sample (TD-learning)    — In-sample (SARSA)

**High dot products**

Breakout

Seaquest

Breakout

**Similar Q-values**

Breakout

Seaquest

Breakout

Incorrect Q-values found eventually; dot-products increase throughout training