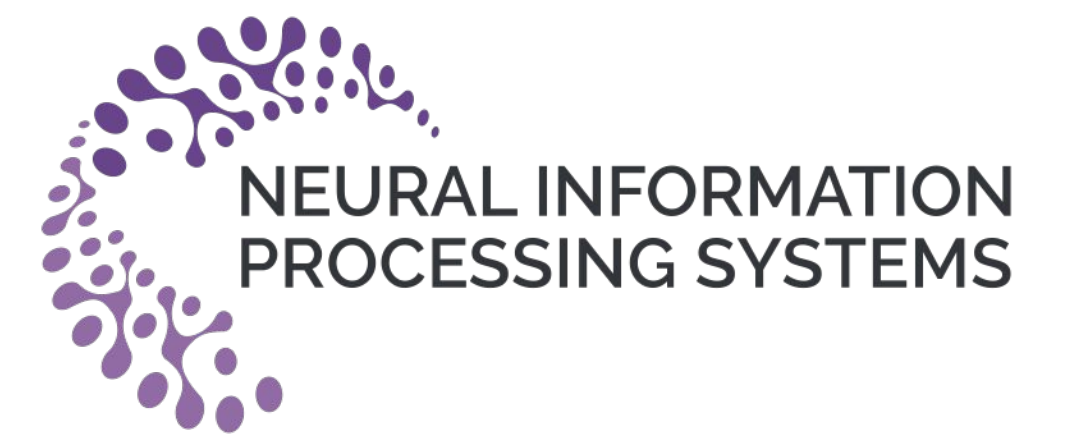


# Count-Based Temperature Scheduling for Maximum Entropy Reinforcement Learning

Dailin Hu,  
Pieter Abbeel,  
Roy Fox



## Motivation

Empirical evidence from Soft Q Learning(SQL)[1] suggests that a state-independent linear scheduling can achieve good performance[2][3].

More insight can be gained from comparing two families of successful RL algorithms:

- G-Learning[2], SQL, Path Consistency Learning(PCL)[4], Soft Actor Critic[5]

$$\pi_i(a|s) \propto \pi_0(a|s) \exp \beta_i(s) Q_i(s, a) \quad \forall s \in \mathcal{S}, a \in \mathcal{A},$$

- Relative Entropy Policy Search[6], Trust Region Policy Optimization[7],

$$\pi_i(a|s) \propto \pi_{i-1}(a|s) \exp \kappa_i(s) Q_i(s, a) \propto \pi_0(a|s) \exp \left( \sum_{j \leq i} \kappa_j(s) Q_j(s, a) \right).$$

Combining the above two equations, we have:

$$\beta_i(s) \approx \sum_{j \leq i} \kappa_j(s) \approx \kappa i,$$

## Count-Based Soft Q Learning

We propose Count-Based Soft Q Learning based on SQL that uses a state-dependent temperature schedule in which  $\beta$  grows linearly with the number of times that the algorithm updates the Q function, for any action.

Let  $n(s,a)$  be the count of sampled data points, then the inverse temperature in CBSQL is

$$\beta(s) = \kappa \sum_a n(s, a)$$

with  $\kappa > 0$  a constant hyperparameter.

## Experiments

### Tabular Experiments

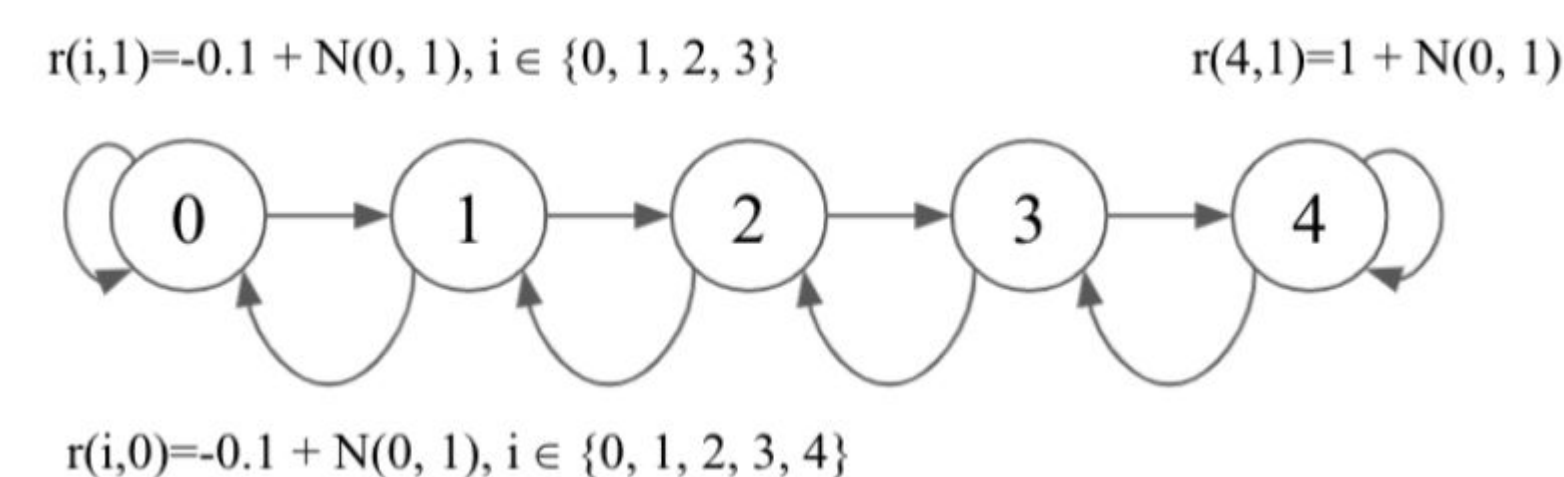


Fig1 . Noisy Chain-walk Problem

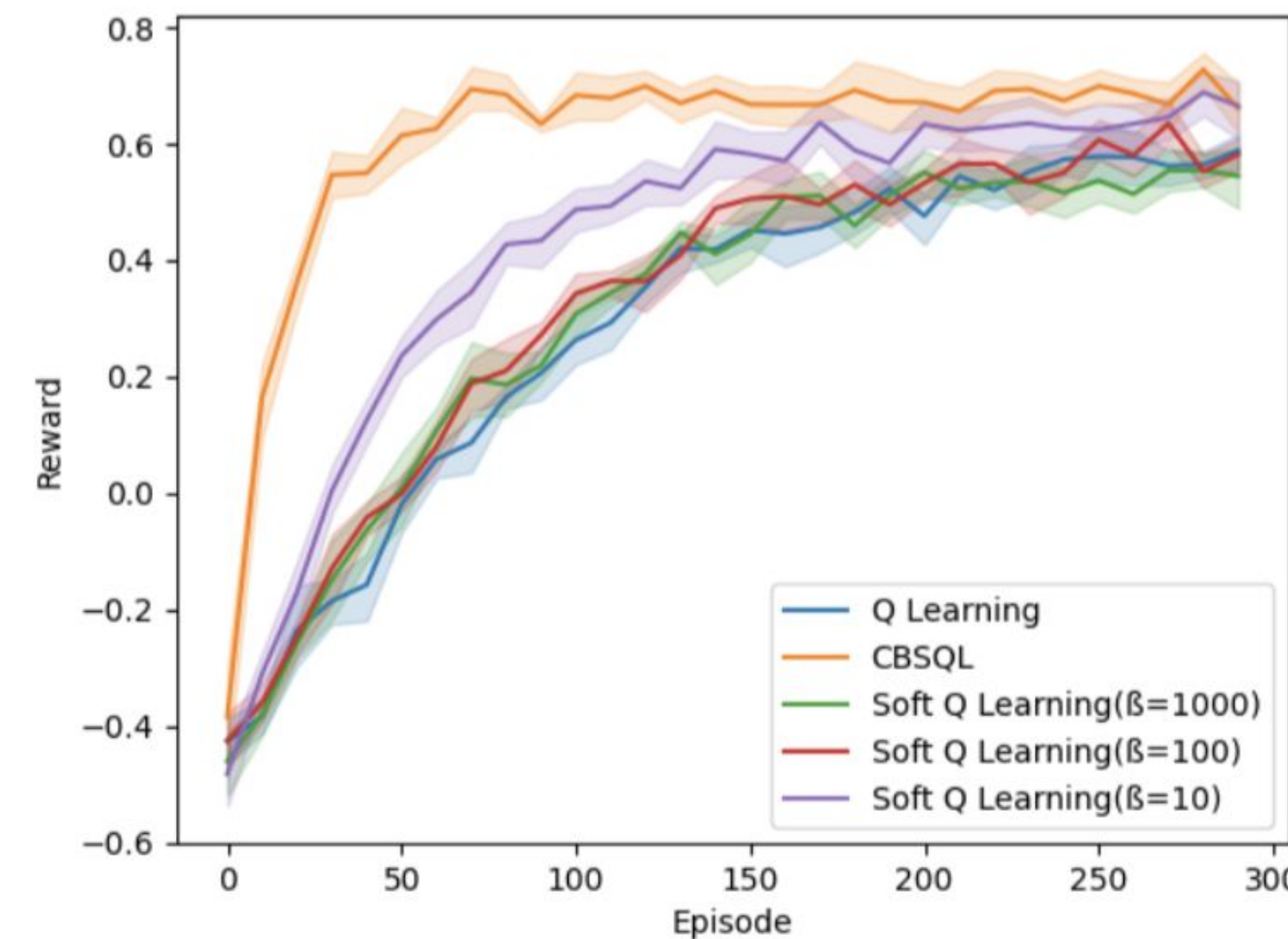


Fig2. Rewards averaged over 10000 runs on the Noisy Chain-Walk Problem

### Atari Experiments

Game	DQN	SQL( $\beta = 100$ )	SQL( $\beta = 1000$ )	CBSQL
Breakout	5.9 ( $\pm 5.9$ )	5.9 ( $\pm 4.5$ )	5.1 ( $\pm 4.7$ )	<b>8.2</b> ( $\pm 6.1$ )
Freeway	21.0 ( $\pm 1.5$ )	14.6 ( $\pm 8.5$ )	22.56 ( $\pm 4.7$ )	<b>25.82</b> ( $\pm 4.9$ )
Pong	1.93 ( $\pm 2.6$ )	<b>17.83</b> ( $\pm 2.2$ )	16.31 ( $\pm 2.7$ )	17.56 ( $\pm 2.0$ )
Qbert	568.4 ( $\pm 1101.9$ )	828 ( $\pm 1411.7$ )	564.5 ( $\pm 1097.5$ )	<b>875.3</b> ( $\pm 1254.6$ )
Seaquest	13.5 ( $\pm 24.1$ )	4 ( $\pm 60.2$ )	17.2 ( $\pm 24.0$ )	<b>84.6</b> ( $\pm 60.2$ )
SpaceInvaders	132.7 ( $\pm 113.2$ )	<b>158.9</b> ( $\pm 128.5$ )	132.25 ( $\pm 118.4$ )	138.9 ( $\pm 112.8$ )

Fig3. DQN, fixed-temperature SQL and CBSQL average rewards (with standard deviation). Raw score are averaged over the last 100 testing episodes across 3 runs.

## Rainbow Integrations to CBSQL

We integrate CBSQL with Rainbow DQN[8], a state-of-the-art reinforcement learning algorithm for memoryless agents including multi-step learning, double-Q learning, prioritized experience replay, dueling networks, distribution RL and noisy networks. All these methods can be straightforwardly applied to soft Q learning except multi-step learning and distributional RL.

### Multi-step learning

Multi-step learning with a tuned-number of steps can lead to faster learning in on-policy RL algorithms. In SQL the  $n$  step truncated return is

$$\tilde{r}_t^{(n)} = r_t^{(n)} + \frac{1}{\beta} \sum_{k=1}^{n-1} \gamma^k \mathbb{H}[\pi(\cdot | s_{t+k})].$$

Unfortunately empirical policy entropy estimates are often very noisy and calls for further study. In this work we simply use 1-step returns for SQL and CBSQL.

### Distributional RL

We adapt distributional RL to SQL and CBSQL by defining a policy distribution of

$$\pi(a'|s') = \frac{\exp \beta(s') \mathbf{z}^\top \mathbf{p}_{\bar{\theta}}(s', a')}{\sum_{\bar{a}'} \exp \beta(s') \mathbf{z}^\top \mathbf{p}_{\bar{\theta}}(s', \bar{a}')}$$

over the values

$$r + \gamma \left( \mathbf{z} - \frac{1}{\beta} \mathbb{D}[\pi(\cdot | s')] \parallel \pi_0 \right)$$

### Preliminary Results

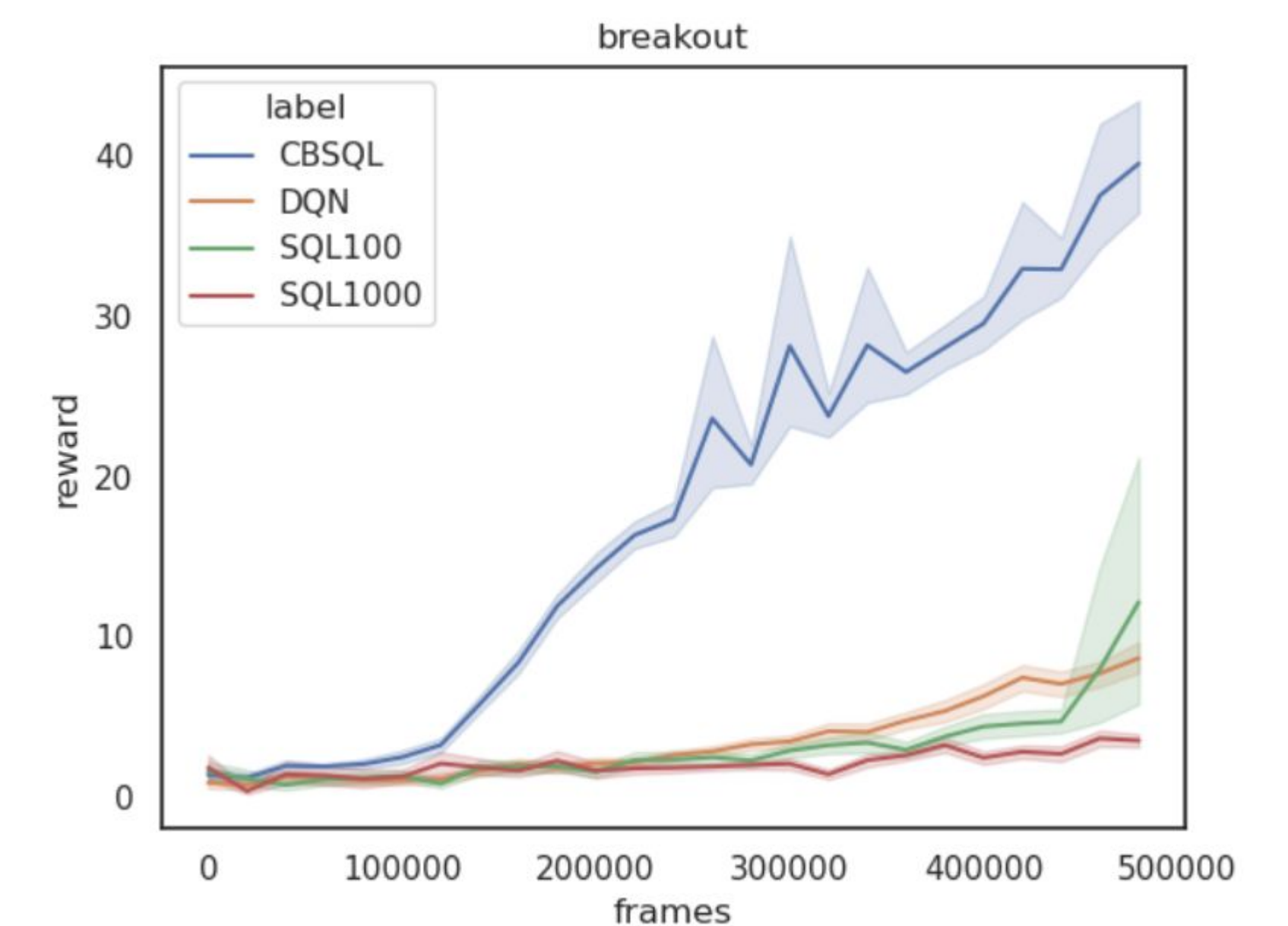


Fig4. CBSQL results compared with DQN and fixed-temperature SQL, with Rainbow. Rewards are averaged over 5 runs.

### References

- [1] Haarnoja T, Tang H, Abbeel P, Levine S. Reinforcement learning with deep energy-based policies. In International Conference on Machine Learning 2017 Jul 17 (pp. 1352-1361). PMLR.
- [2] Fox R, Pakman A, Tishby N. Taming the noise in reinforcement learning via soft updates. arXiv preprint arXiv:1512.08562. 2016
- [3] Grau-Moya J, Leibfried F, Vrancx P. Soft q-learning with mutual-information regularization. In International conference on learning representations 2018 Sep 27.
- [4] Chow Y, Nachum O, Ghavamzadeh M. Path consistency learning in tsallis entropy regularized mdps. In International Conference on Machine Learning 2018 Jul 3 (pp. 979-988). PMLR.
- [5] Haarnoja T, Zhou A, Abbeel P, Levine S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In International conference on machine learning 2018 Jul 3 (pp. 1861-1870). PMLR.
- [6] Peters J, Mulling K, Altun Y. Relative entropy policy search. In Twenty-Fourth AAAI Conference on Artificial Intelligence 2010 Jul 5.
- [7] Schulman J, Levine S, Abbeel P, Jordan M, Moritz P. Trust region policy optimization. In International conference on machine learning 2015 Jun 1 (pp. 1889-1897). PMLR.
- [8] Hessel M, Modayil J, Van Hasselt H, Schaul T, Ostrovski G, Dabney W, Horgan D, Piot B, Azar M, Silver D. Rainbow: Combining improvements in deep reinforcement learning. In Thirty-second AAAI conference on artificial intelligence 2018 Apr 29.

