

# Preliminary Experiment

January 17, 2018

Bandit Prior/Instances that were considered (only for 10 arms):

- Uniform: Draw the mean rewards for the arms from  $[0.25, 0.75]$
- “HeavyTail”: We took the mean rewards to be randomly drawn from  $\text{Beta}(\alpha = 0.6, \beta = 0.6)$ . With this distribution it was likely to have arms that were at the extremes (close to 1 and close to 0) but also some of the arms with intermediate value means.
- Needle-in-haystack
  1. Low Mean
    - (a) Low - 9 arms with mean 0.01, 1 arm with mean 0.02 (+ 0.01)
    - (b) Medium - 9 arms with mean 0.01, 1 arm with mean 0.06 (+ 0.05)
    - (c) High - 9 arms with mean 0.01, 1 arm with mean 0.21 (+ 0.20)
  2. Medium Mean
    - (a) Low - 9 arms with mean 0.50, 1 arm with mean 0.01 (+ 0.01)
    - (b) Medium - 9 arms with mean 0.50, 1 arm with mean 0.55 (+ 0.05)
    - (c) High - 9 arms with mean 0.50, 1 arm with mean 0.70 (+ 0.20)
  3. High Mean
    - (a) Low - 9 arms with mean 0.80, 1 arm with mean 0.81 (+ 0.01)
    - (b) Medium - 9 arms with mean 0.80, 1 arm with mean 0.85 (+ 0.05)
    - (c) High - 9 arms with mean 0.80, 1 arm with mean 1.0 (+ 0.20)

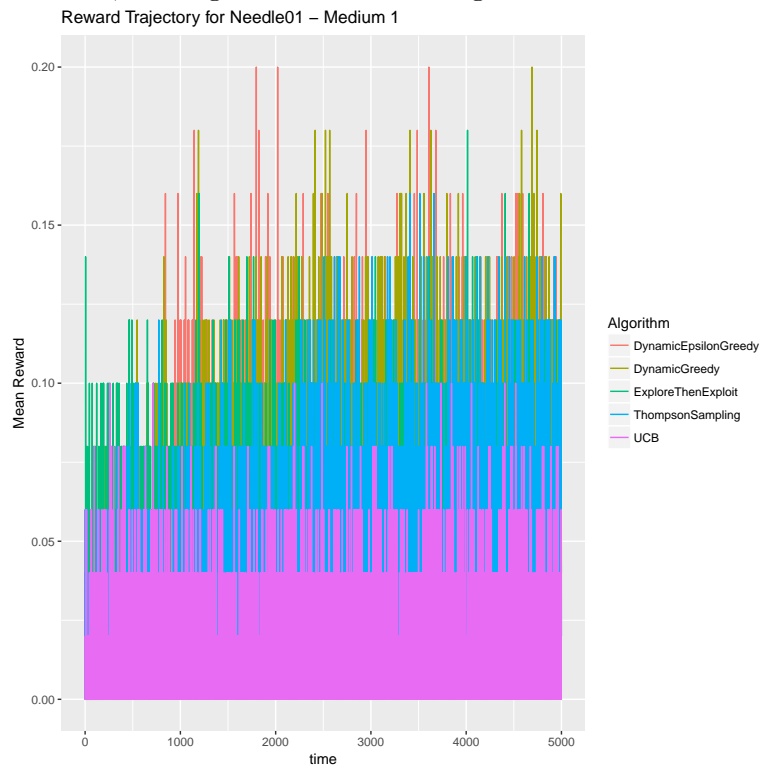
Algorithms considered:

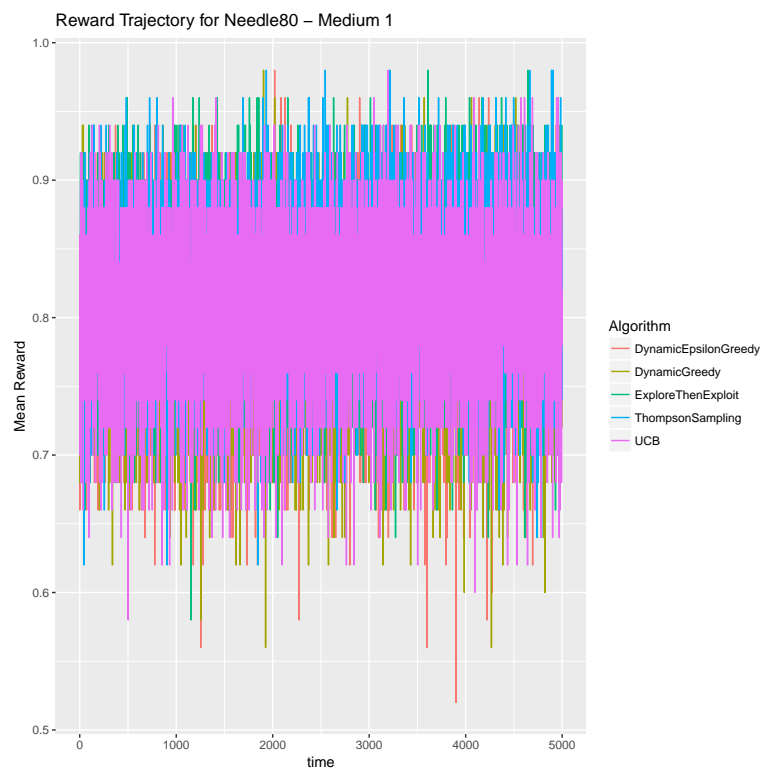
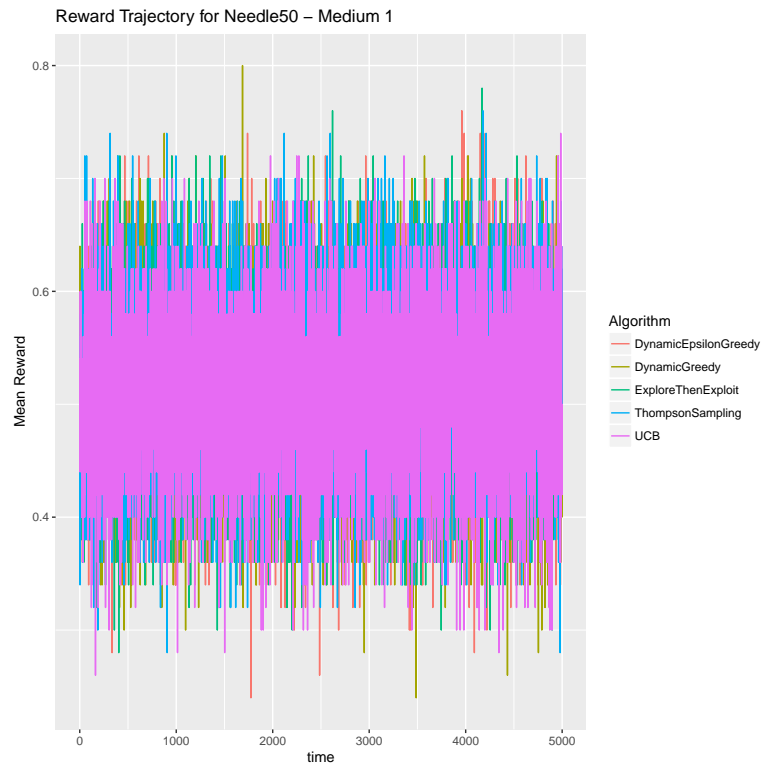
1. ThompsonSampling with priors of  $\text{Beta}(5, 5)$  for every arm.
2. ExploreThenExploit with 100 exploration steps
3. DynamicGreedy with priors of  $\text{Beta}(5, 5)$  for every arm
4. Dynamic  $\epsilon$ -greedy with  $\epsilon=0.05$
5. UCB1

The first run considered  $N = 50$  runs. The main use of this run was to help us detect any qualitative differences between the different needle-in-haystack regimes we Originally we considered three regimes of “needle-in-haystack” instances since we weren’t sure if learning would be different between the different regimes. We considered so many needle-in-haystack instances to see if there would be meaningful qualitative differences between learning when the means were at extremes vs when most means were in the middle of the range of possible means.

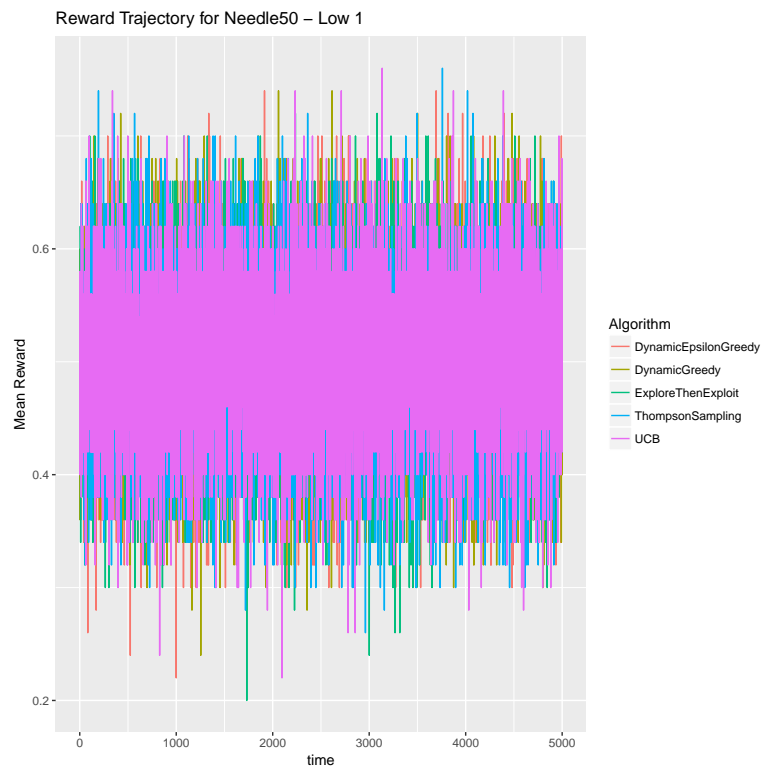
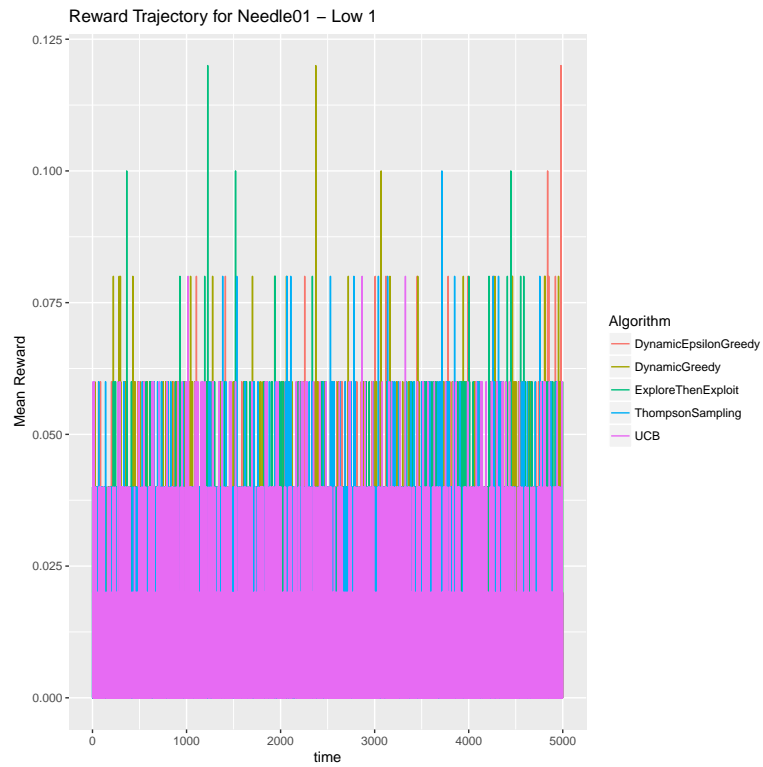
The plots for all the algorithms are reported below. The “low”  $\Delta$  in the needle-in-haystack appears uninteresting across all considered means, likely because the  $\Delta$  of 0.01 is too small to be meaningful. For the “medium” instances it appears as though there is a qualitative difference between the “low” mean and the others. Though this may be partially due to low  $N$ , it appears as though there is some divergence particularly between UCB and Thompson Sampling. For the high  $\Delta$  it the “low” mean and “medium” mean seem to have qualitatively similar results. The “high” mean has starker differences in terms of the learning rates of the algorithms.

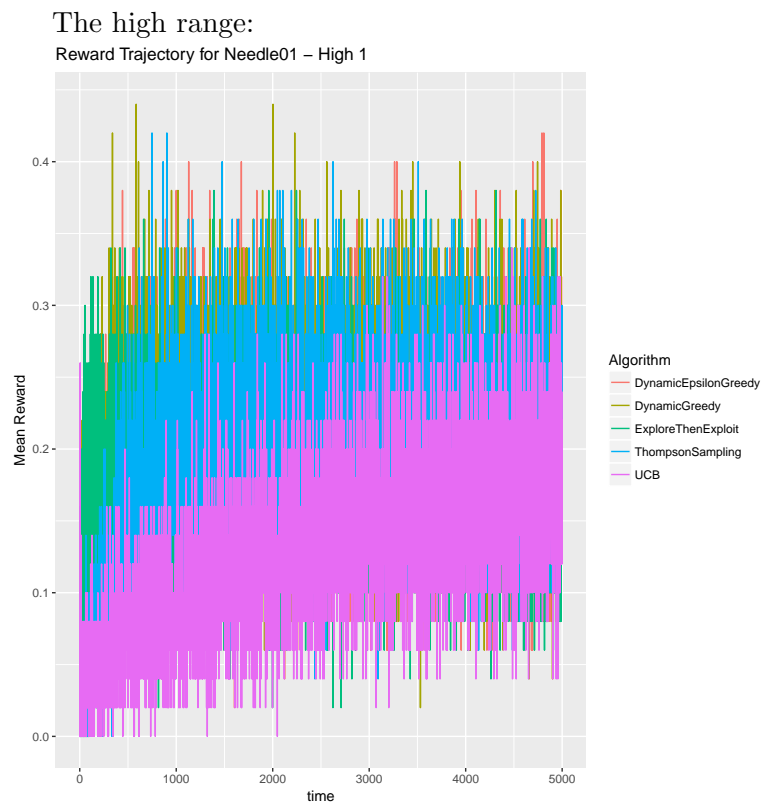
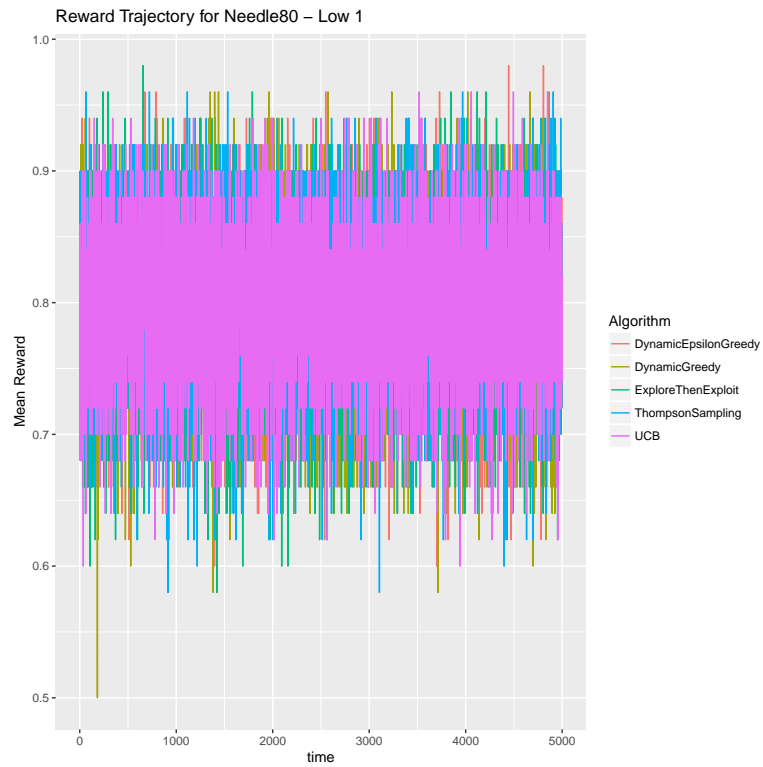
First, looking at the medium range:

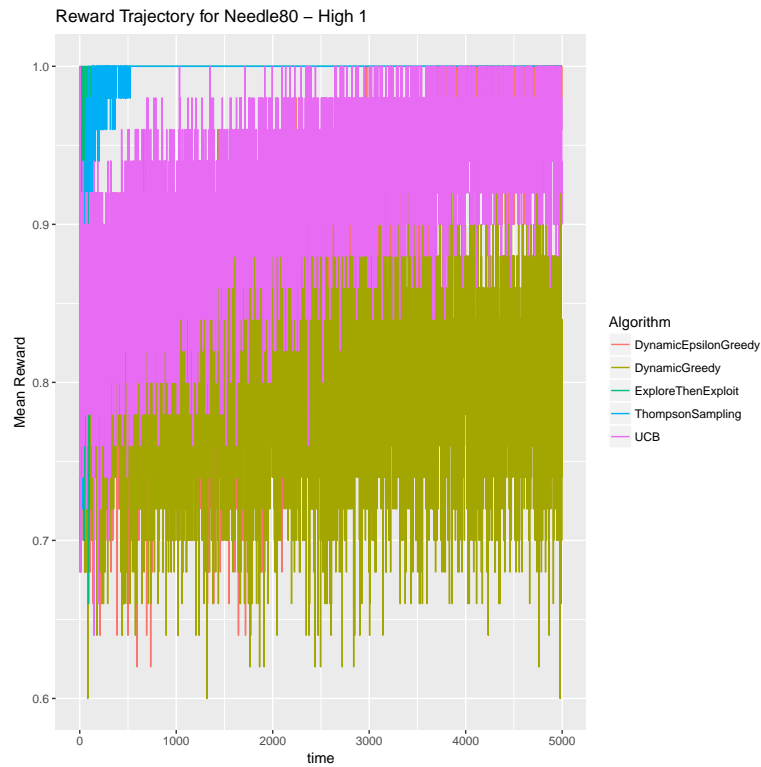
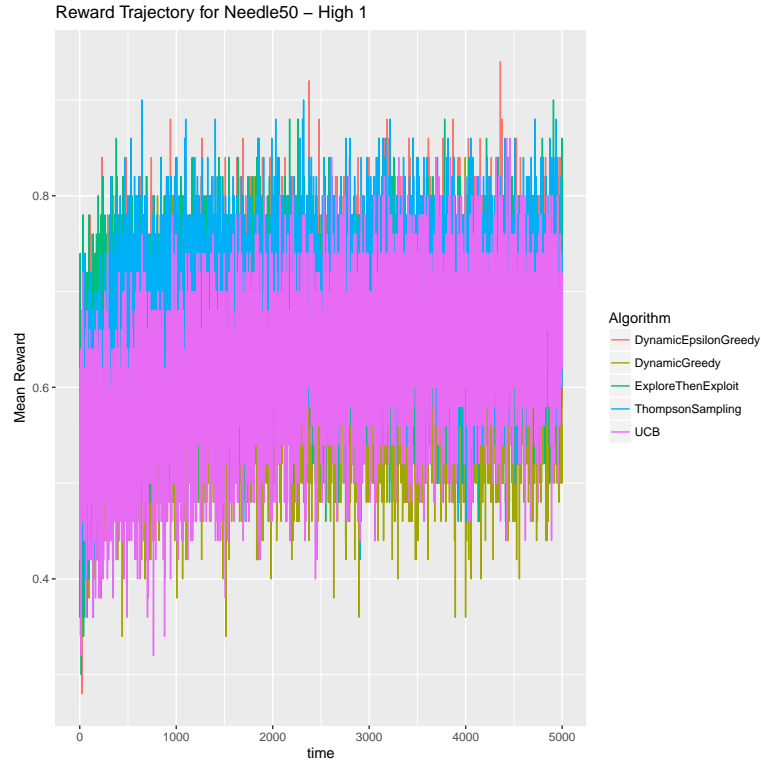




The low range:



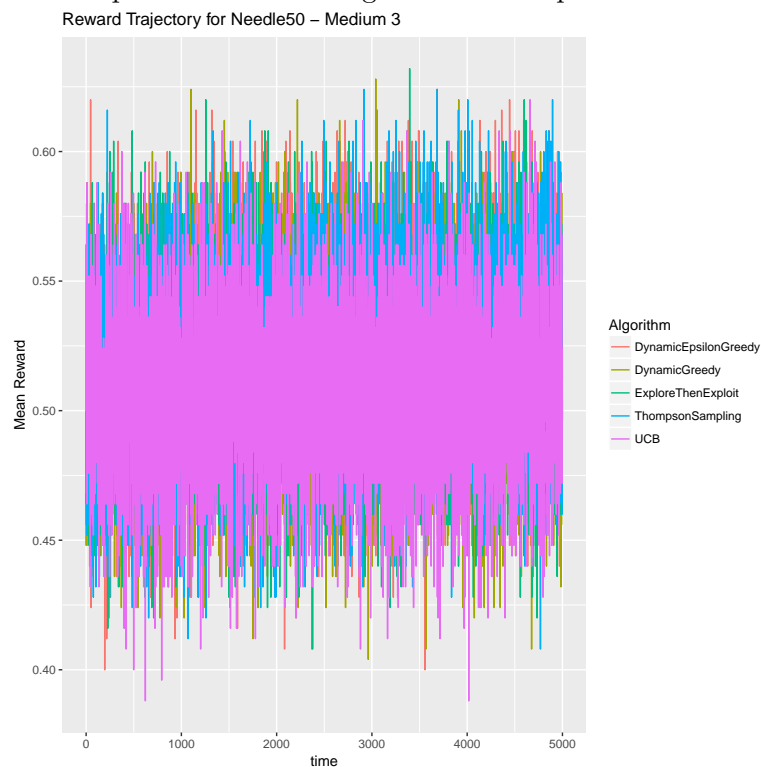


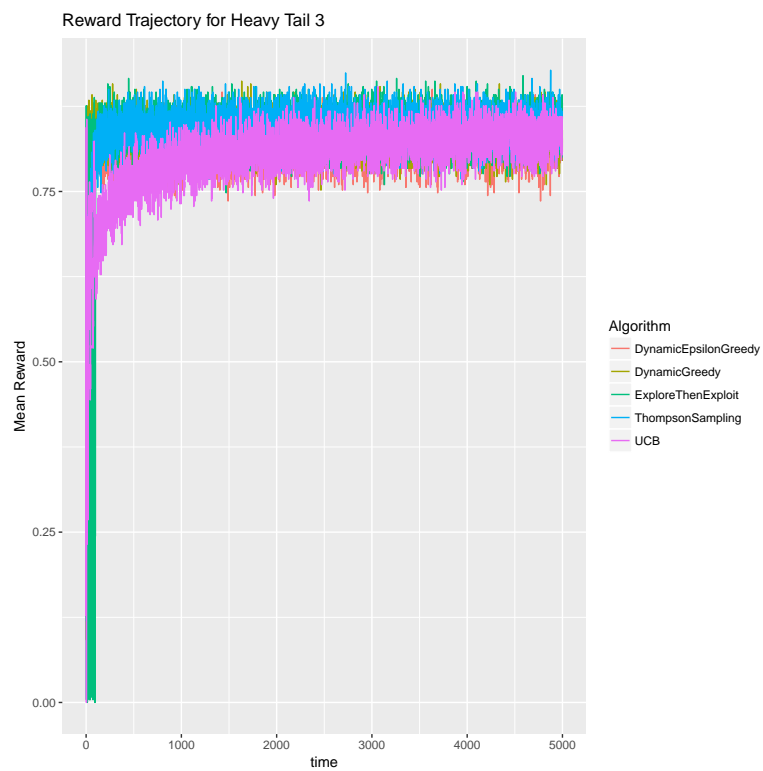
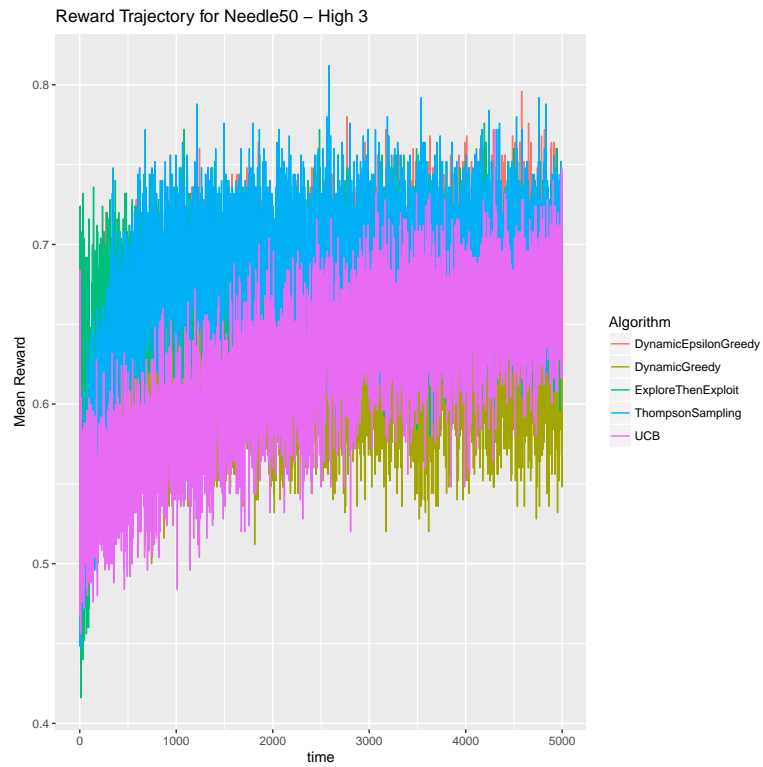


Surprisingly with  $N = 50$  there appeared to be considerable variability in the trajectories and so the experiment was reran with  $N = 250$  and  $T = 5000$ . To reduce the amount of time this

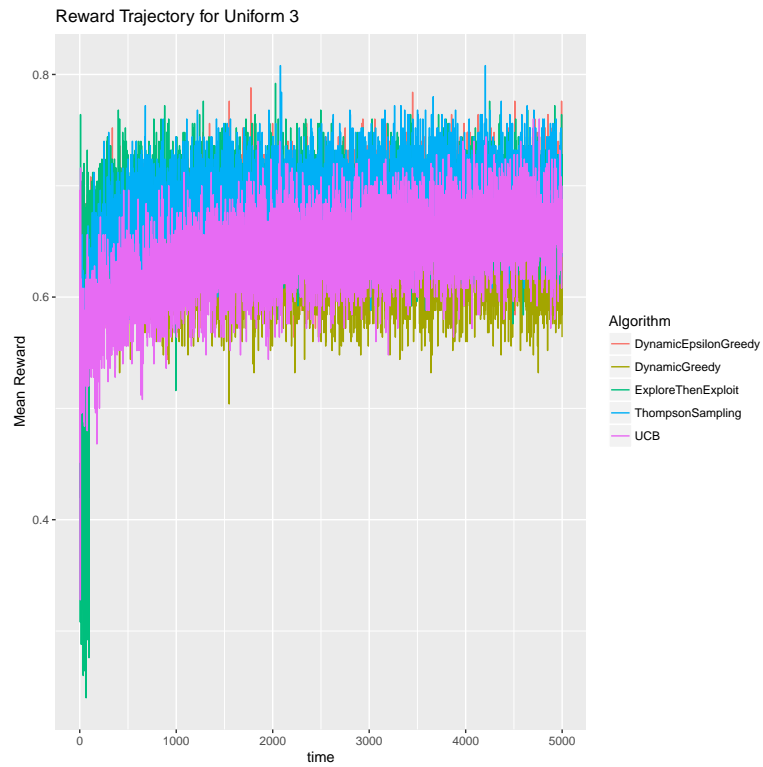
took, only Uniform, HeavyTail, Needle50 - Medium, and Needle50 - High were considered. Though there were some qualitative differences between the different needle in haystack instances discussed above, they were omitted here but we can run this again with those if we decide we want some more fine-grain information.

The plots with all the algorithms are reported below:

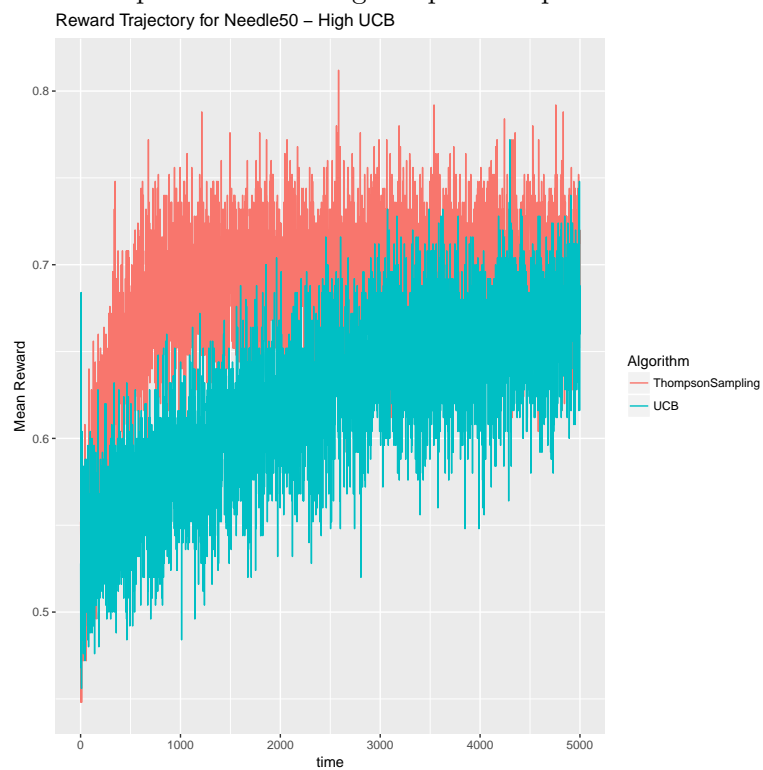


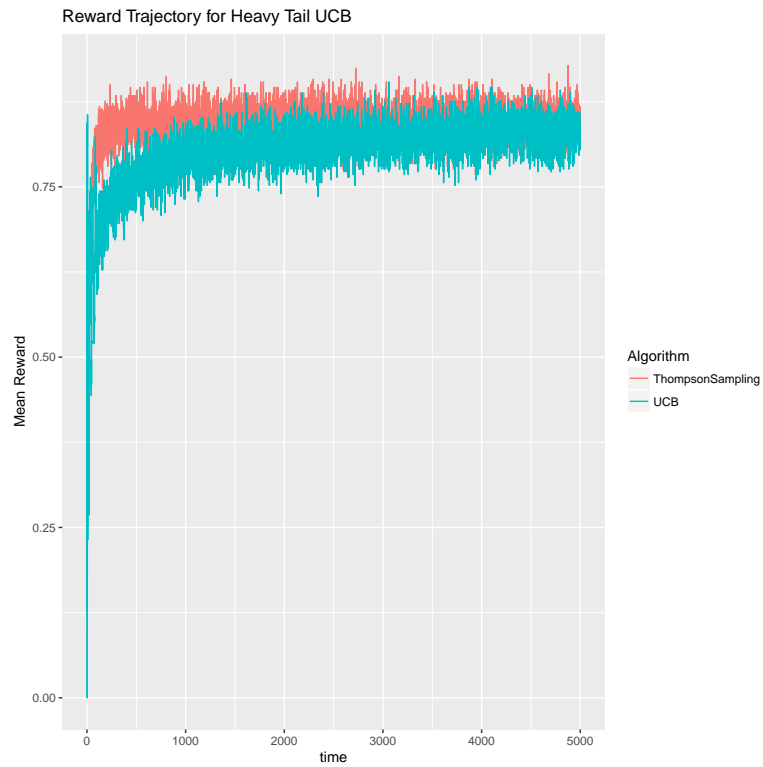




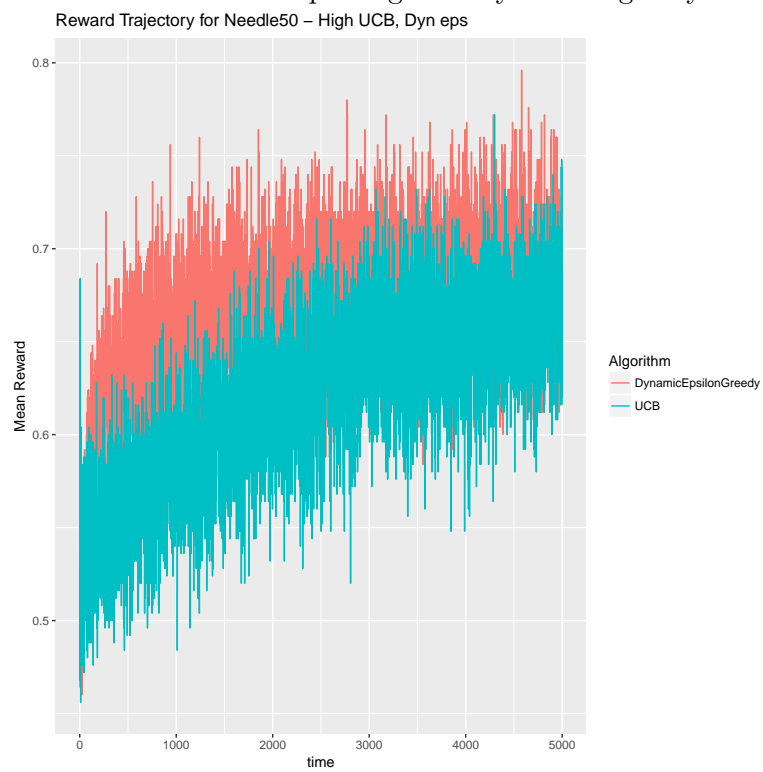


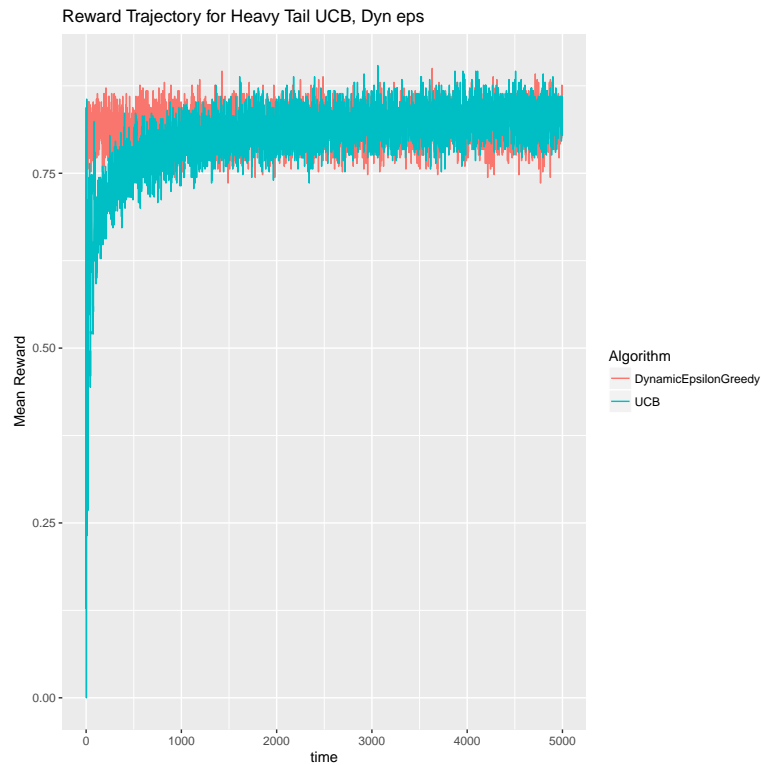
A few particular striking comparative plots:





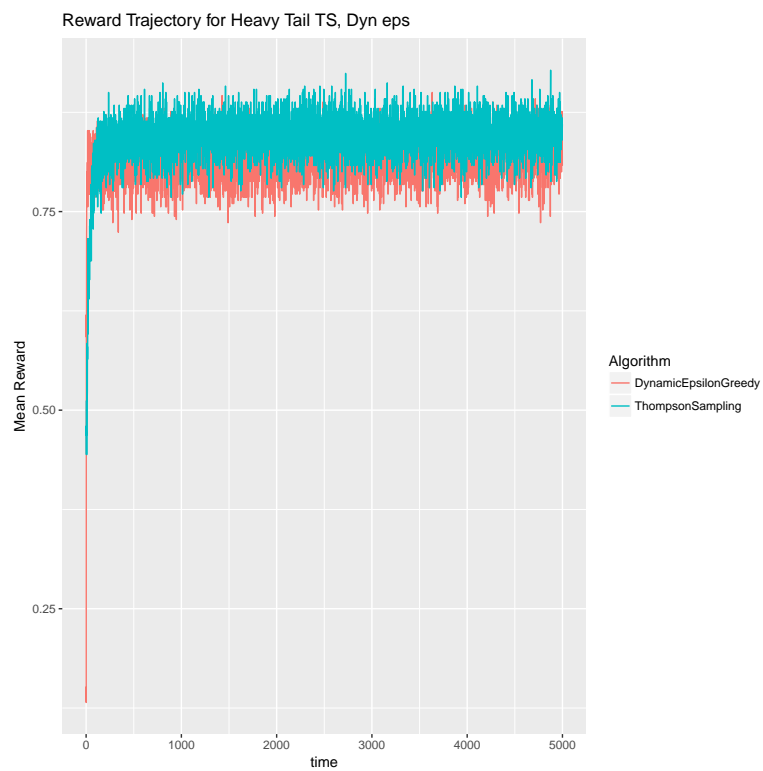
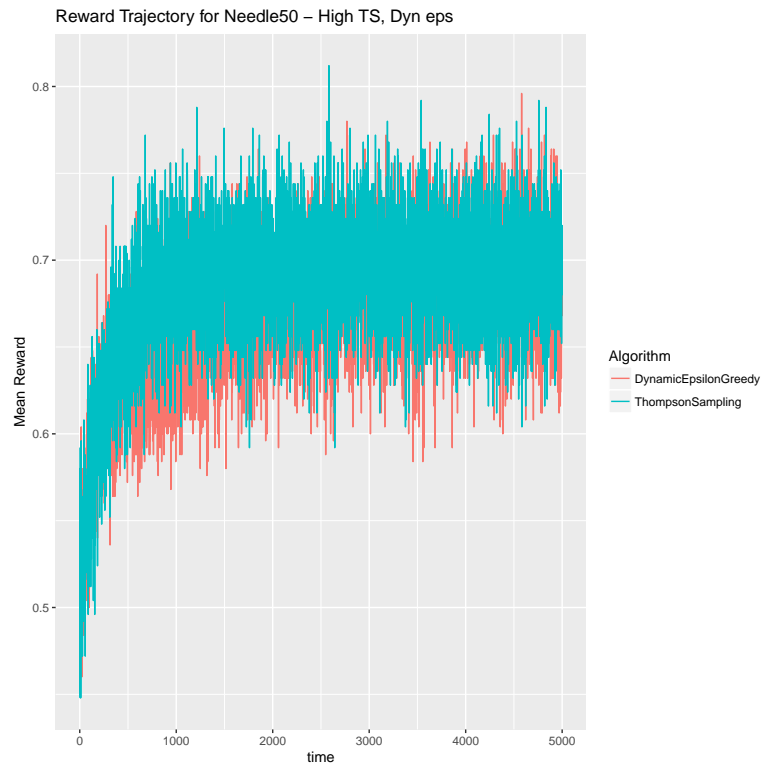
How does UCB compare against Dynamic  $\epsilon$ -greedy in these cases?





The rate of learning for UCB seems slower than for the other algorithms. Is this expected (or do we need to tune the UCB1 constant - I don't think so)? This potentially explains the results we saw that UCB1 loses in the competing bandits game because its rate of learning growth for these bandit instances is slower than the other algorithms.

Comparing ThompsonSampling to Dynamic  $\epsilon$ -greedy:



As well, a version of the plots that smooths the trajectory:

