

Signed Maximum Logit Change; $S_{adv} - S$

0.2
0.0
-0.2
-0.4
-0.6

Familiar

Novel

