

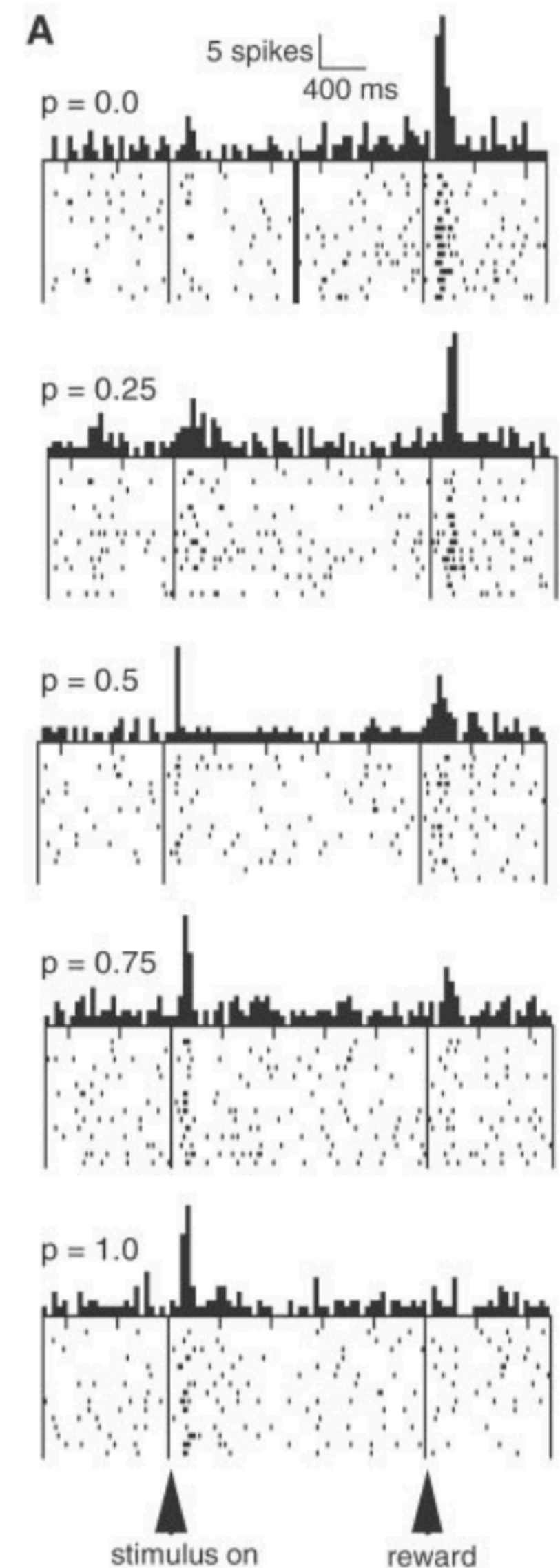
Distributional RL

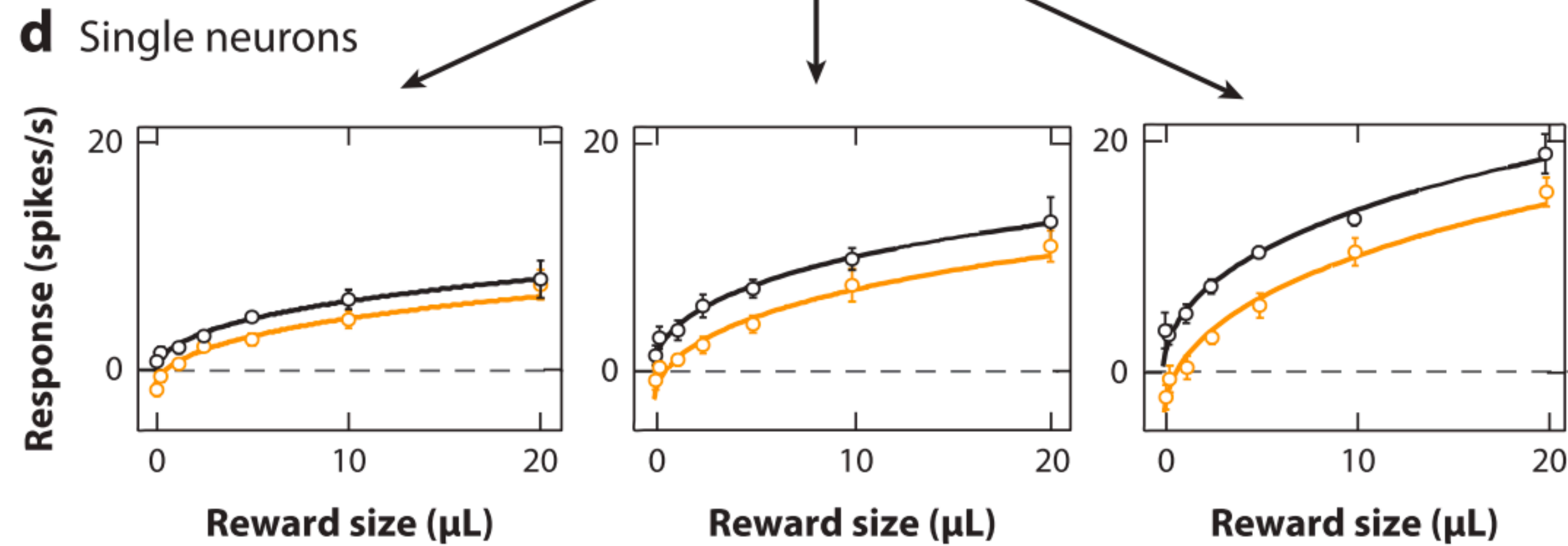
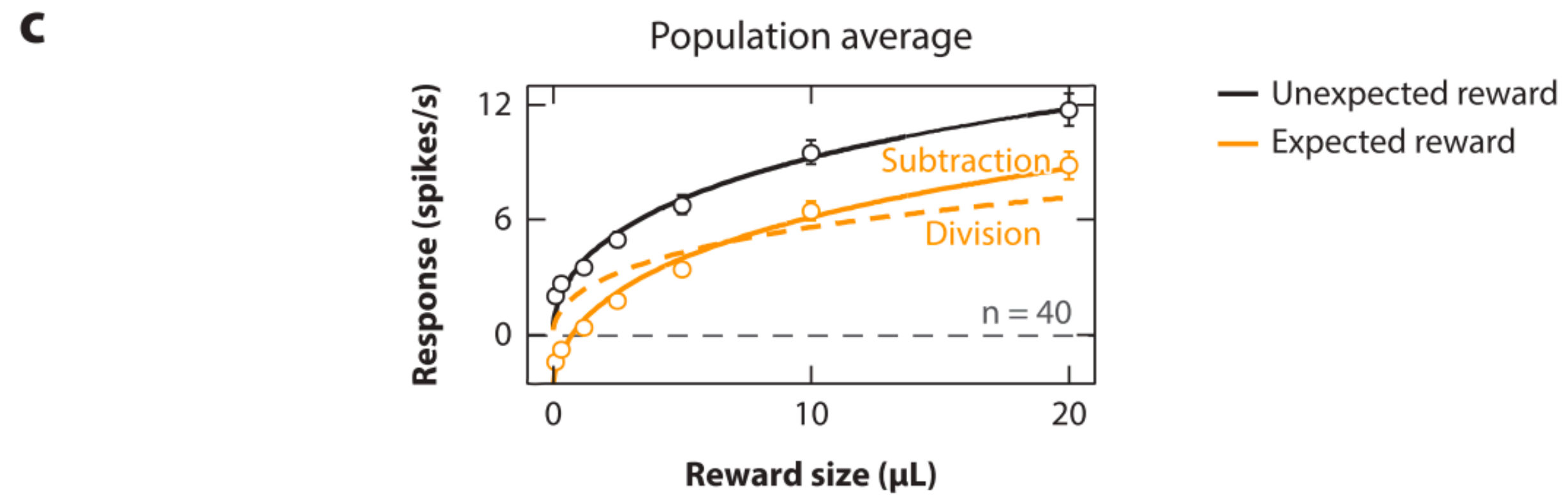
Brabeeba Wang

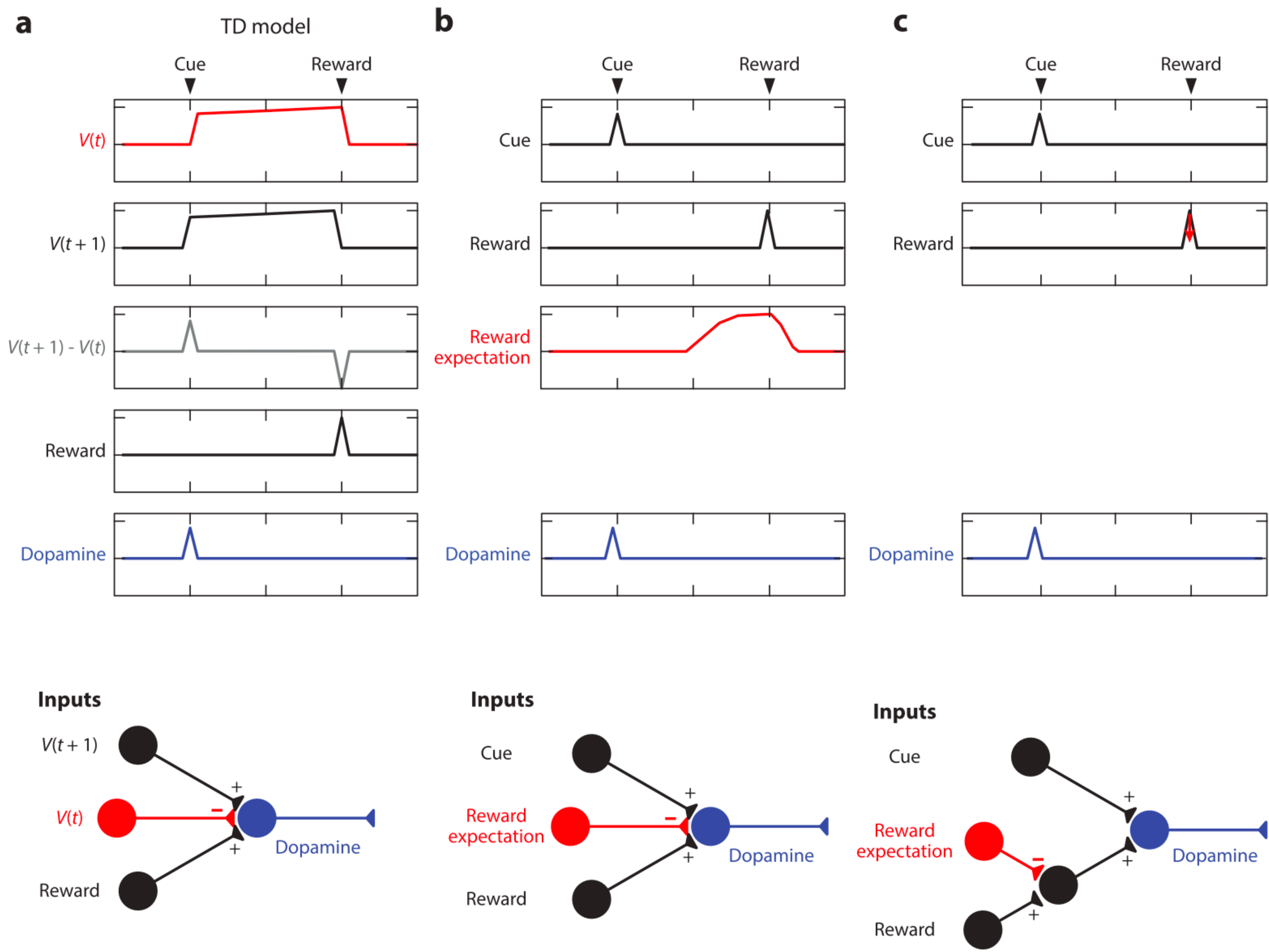
4/9/2021

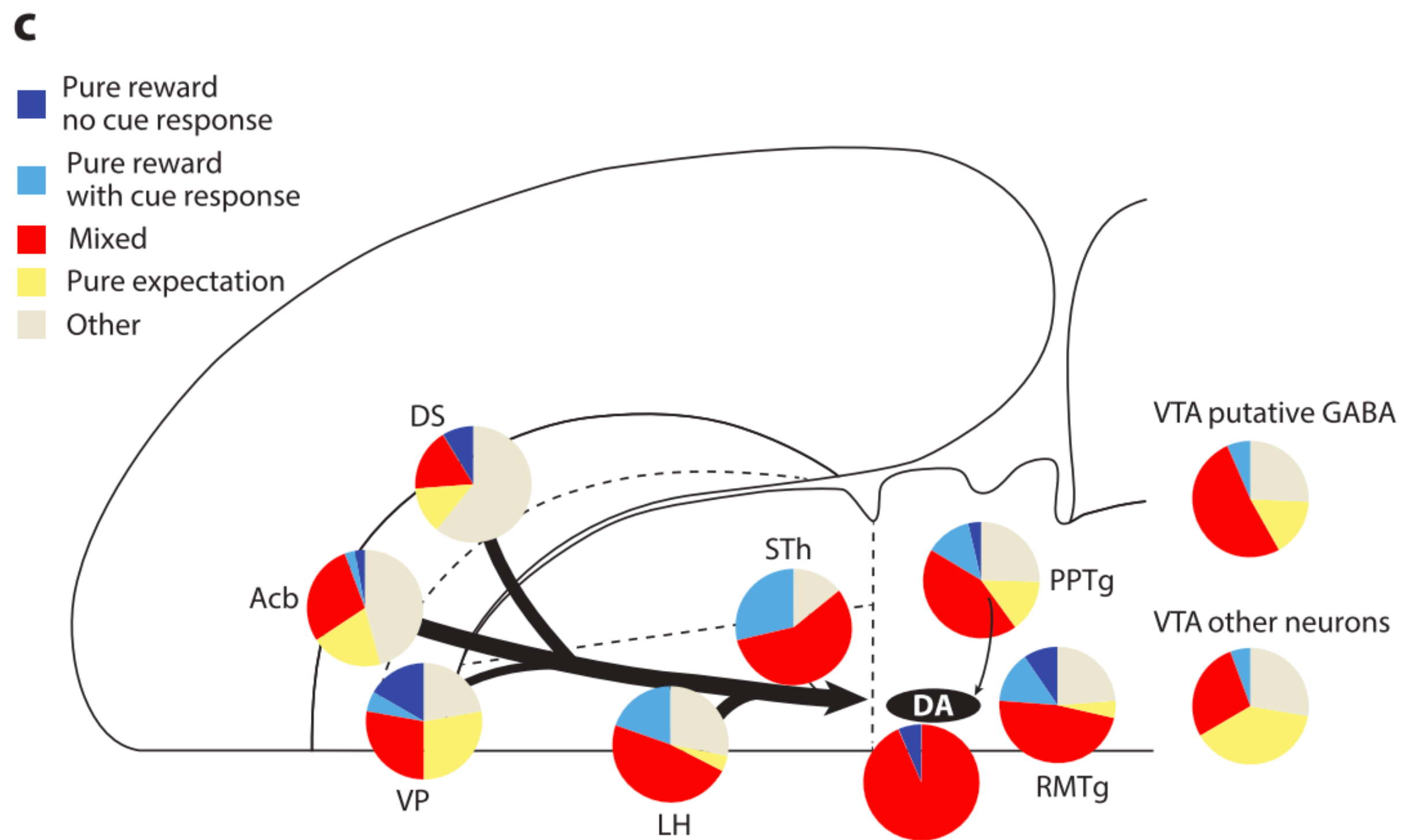
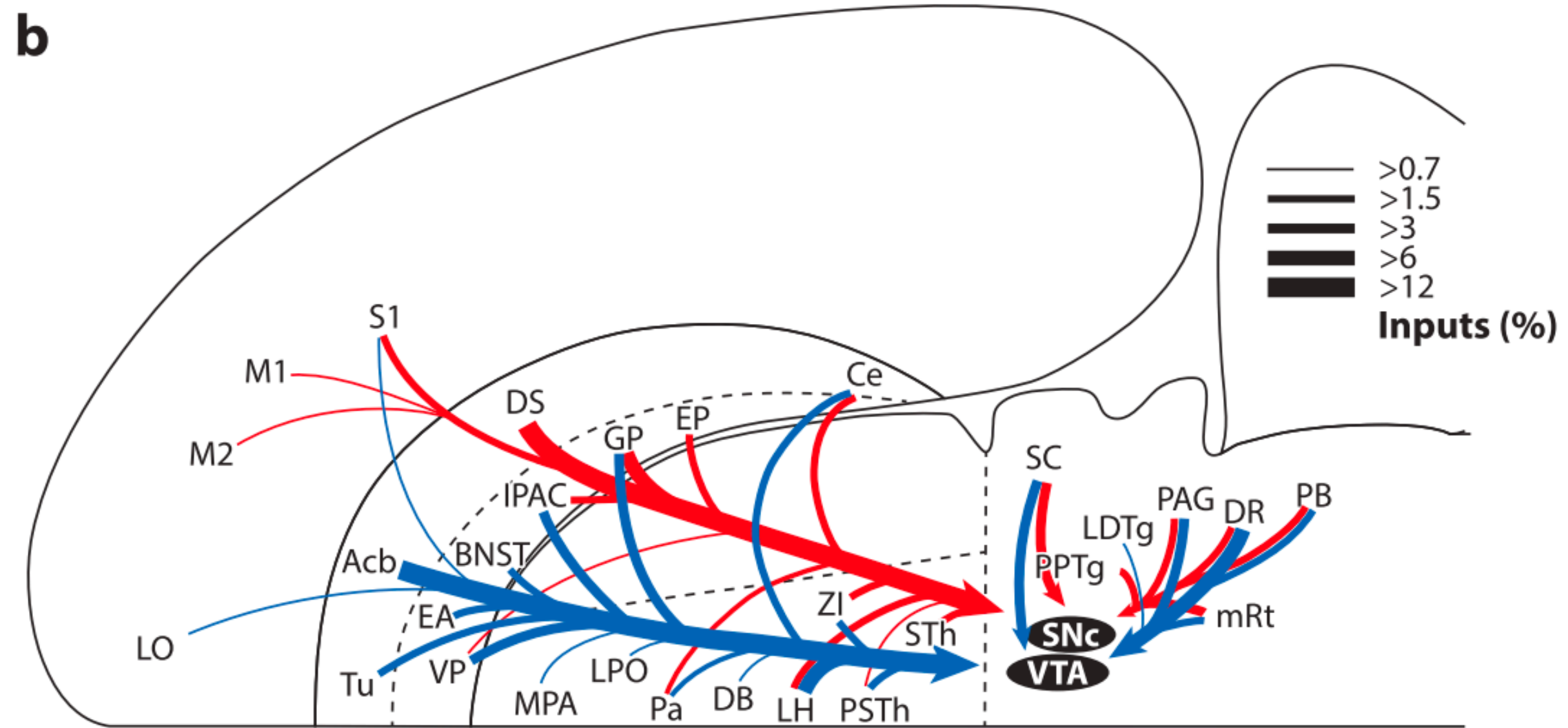
Dopamine

- One of the primary role of dopamine is to encode reward prediction error.
- $\delta = R(s) + \gamma V(s') - V(s)$
- This serves as a reinforcer for learning.
- By optogenetically activate DA neurons, one can reinforce desired behaviors in the animals



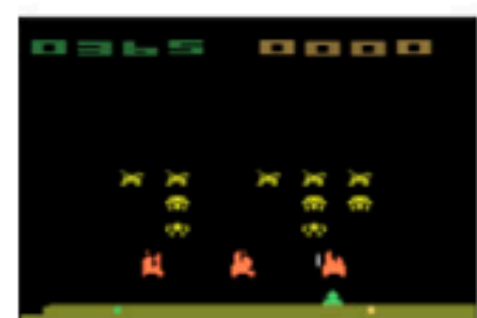






(C)

Traditional RL



State



Expected value
(mean)

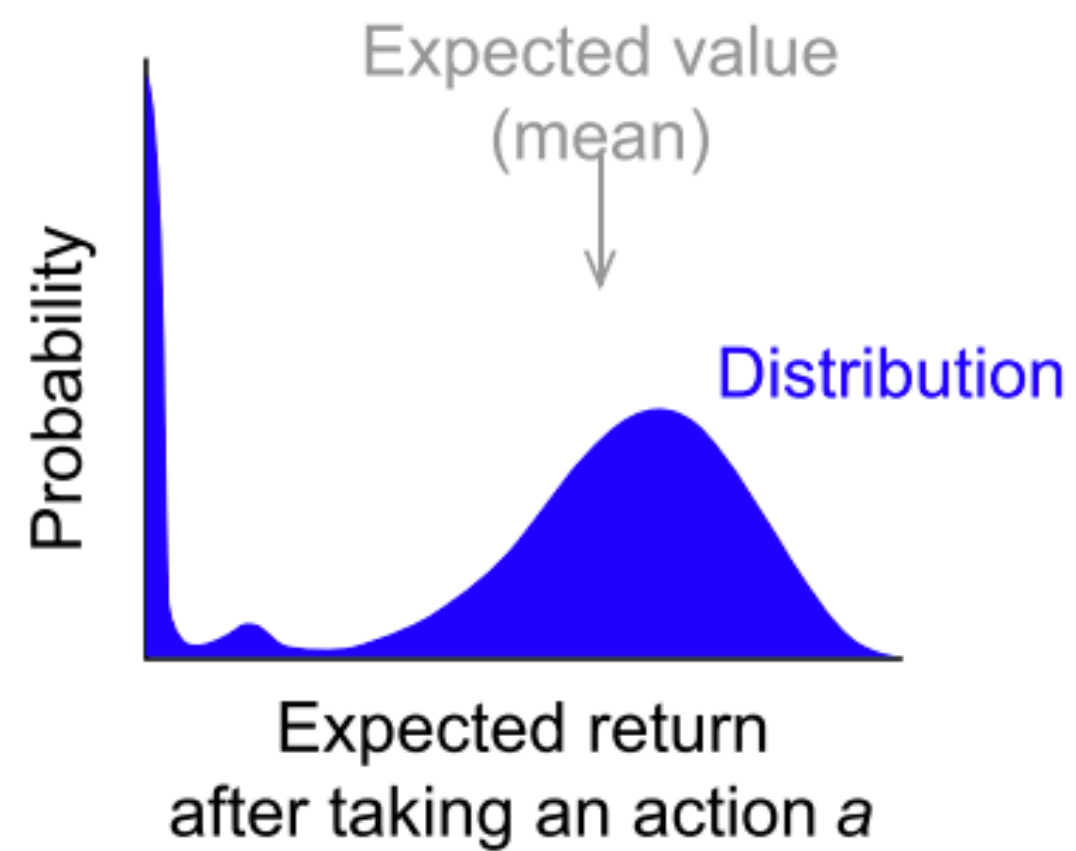
Distributional RL



State

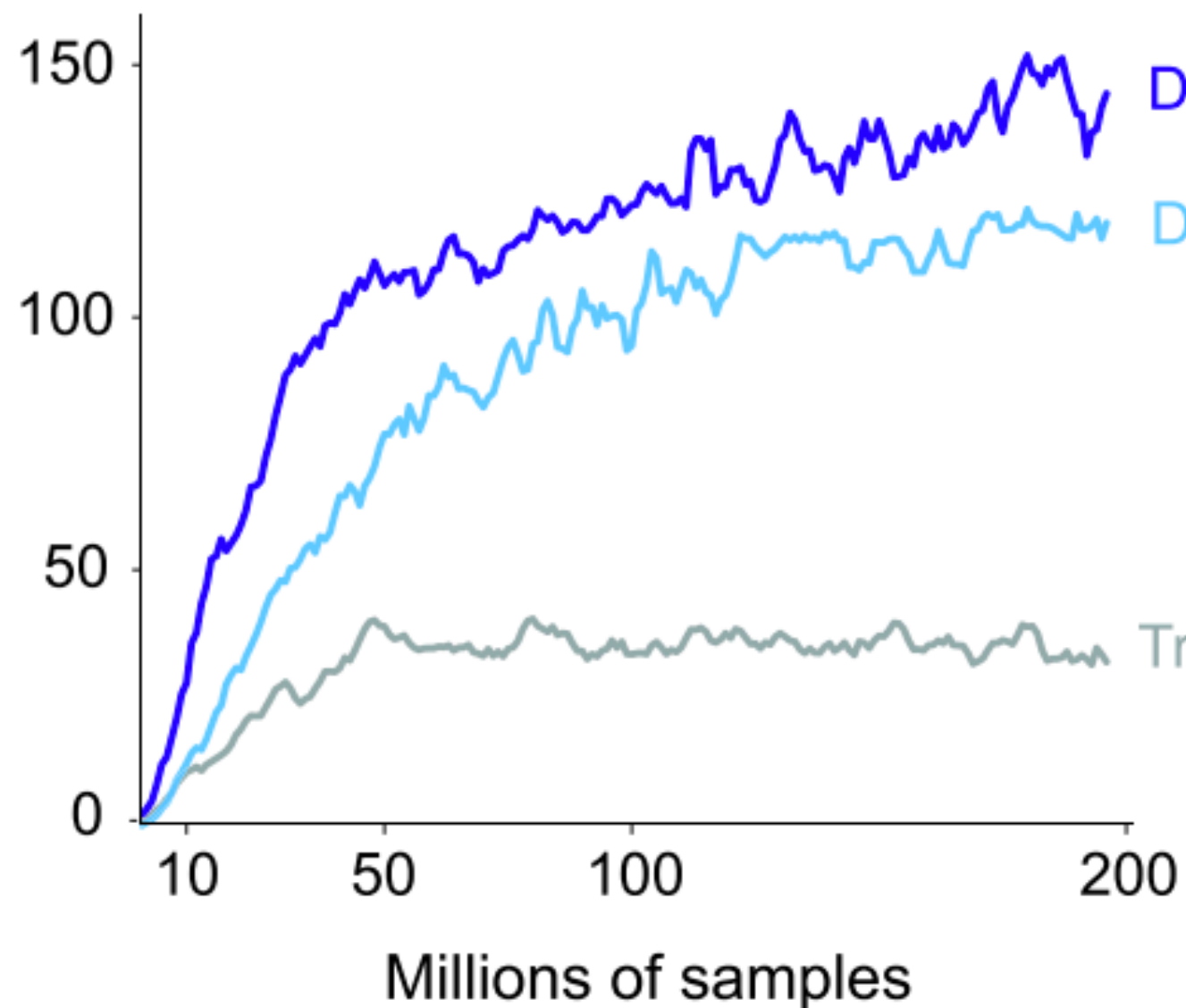


Distribution



(D)

Median human normalized score (%)

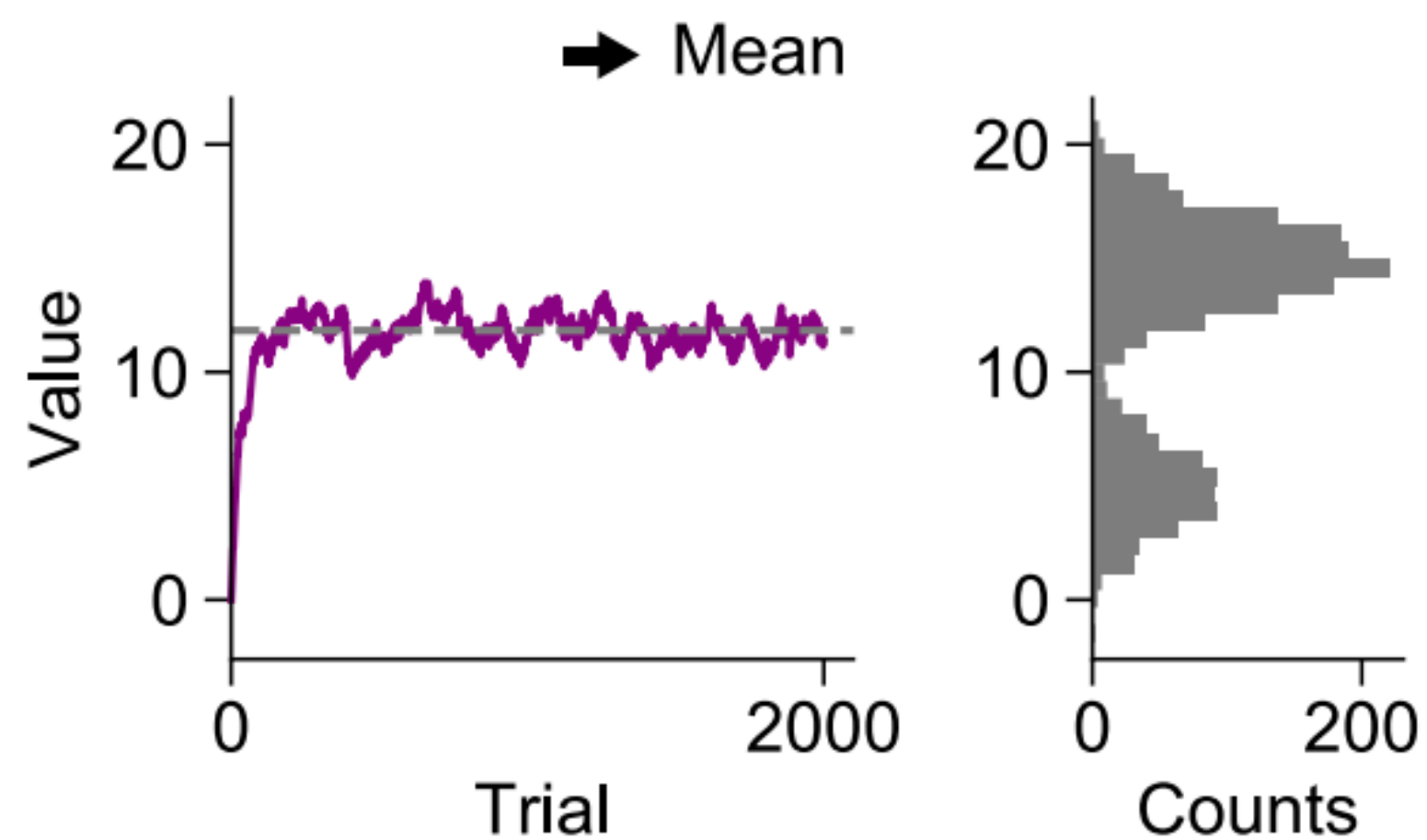


Distributional RL (quantile)

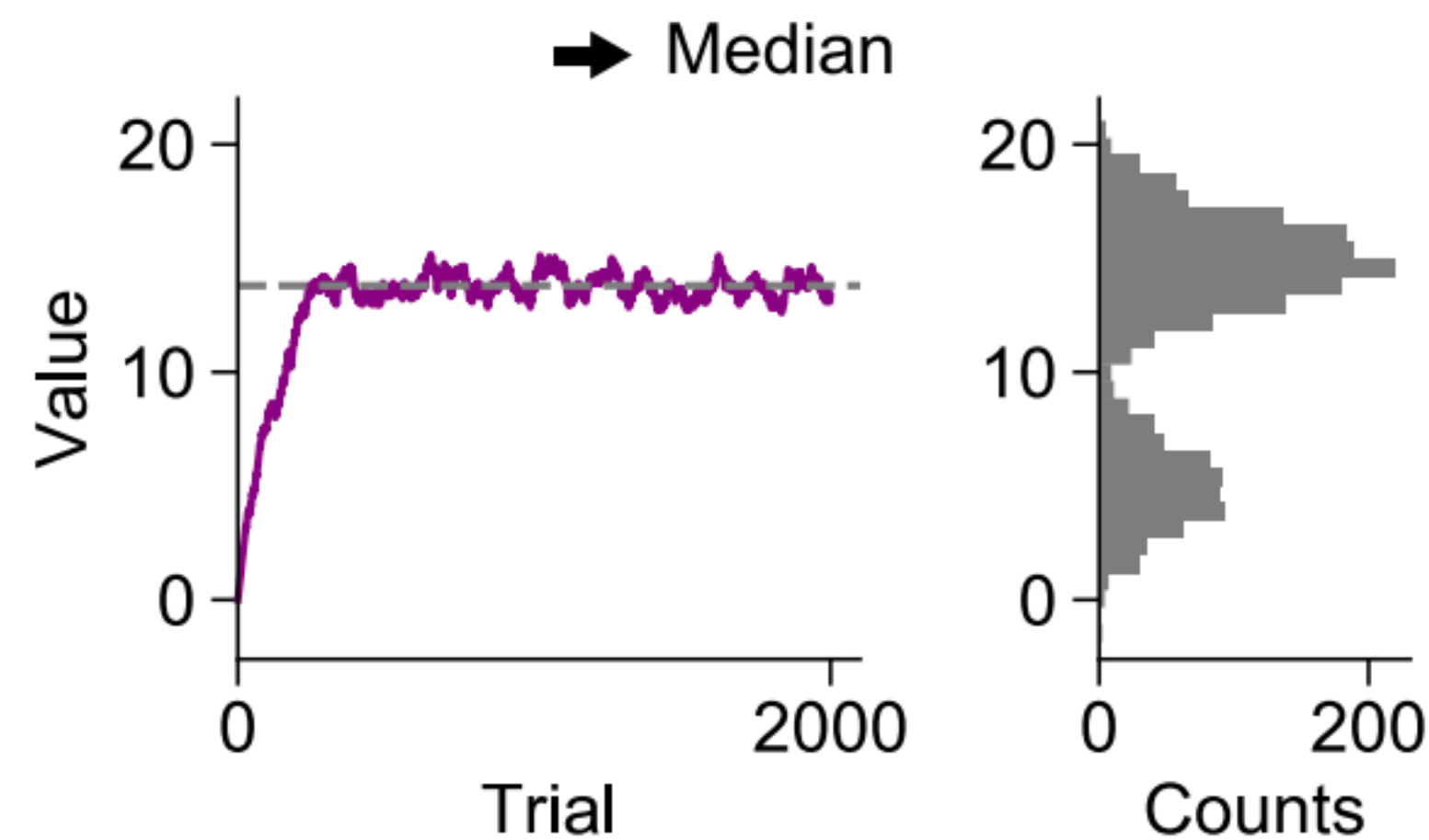
Distributional RL (categorical)

Traditional RL (DQN)

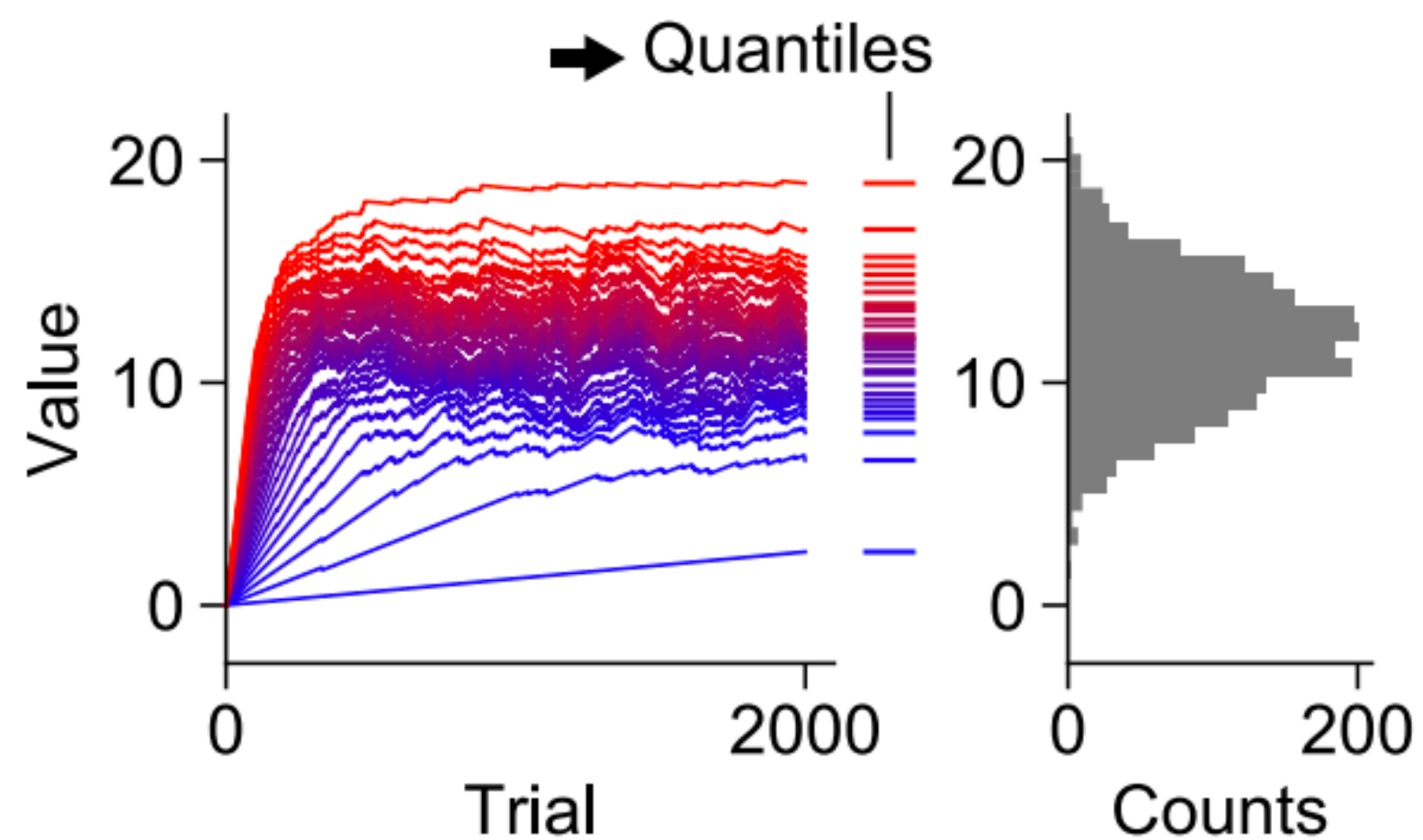
(A) $V \leftarrow V + \alpha \cdot \delta$



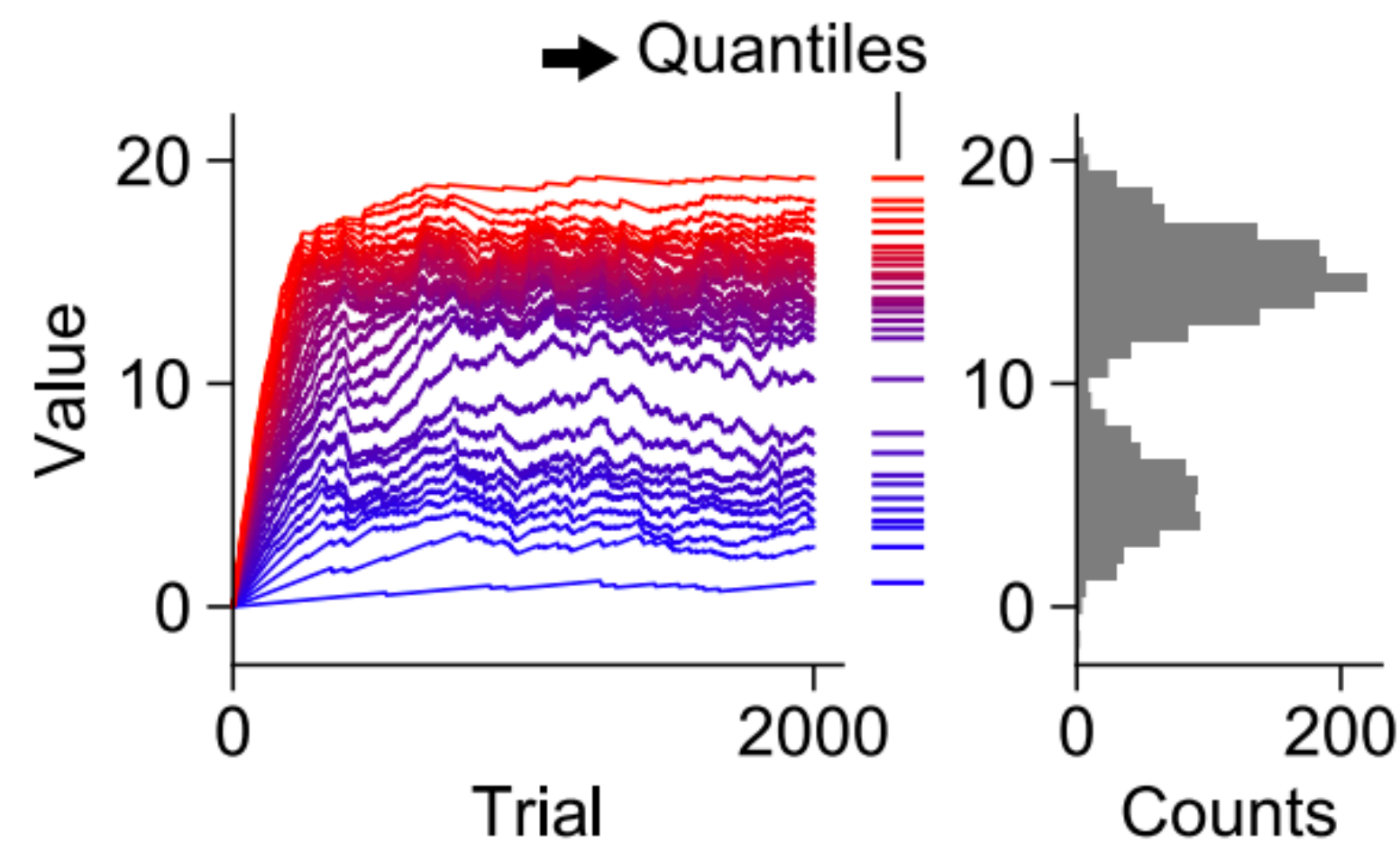
(B) $V \leftarrow V + \alpha \cdot \begin{cases} -1 & \text{if } \delta \leq 0 \\ 1 & \text{if } \delta > 0 \end{cases}$



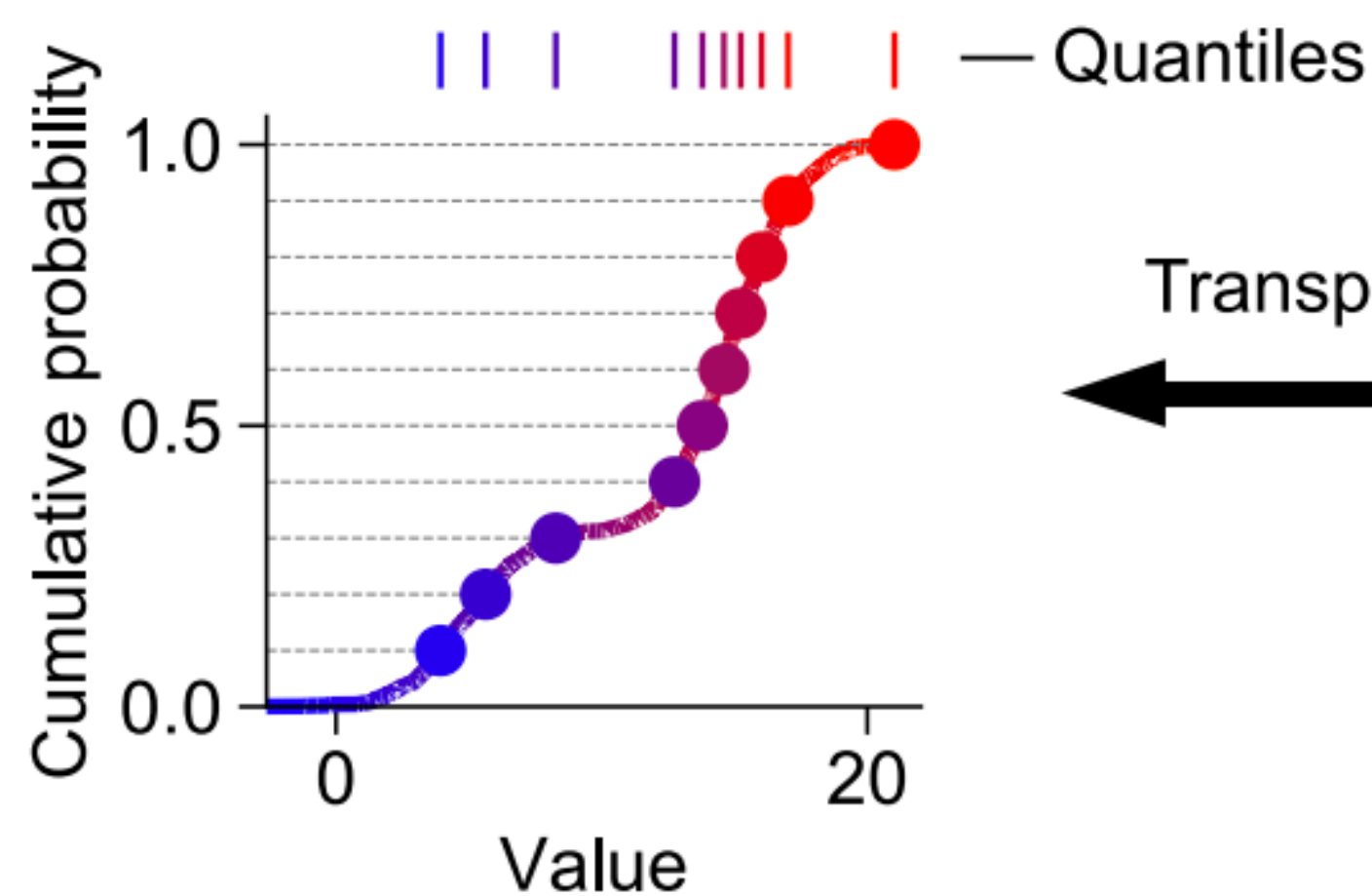
(C) $V_i \leftarrow V_i + \begin{cases} \alpha_i^- \cdot (-1) & \text{if } \delta_i \leq 0 \\ \alpha_i^+ \cdot (+1) & \text{if } \delta_i > 0 \end{cases}$



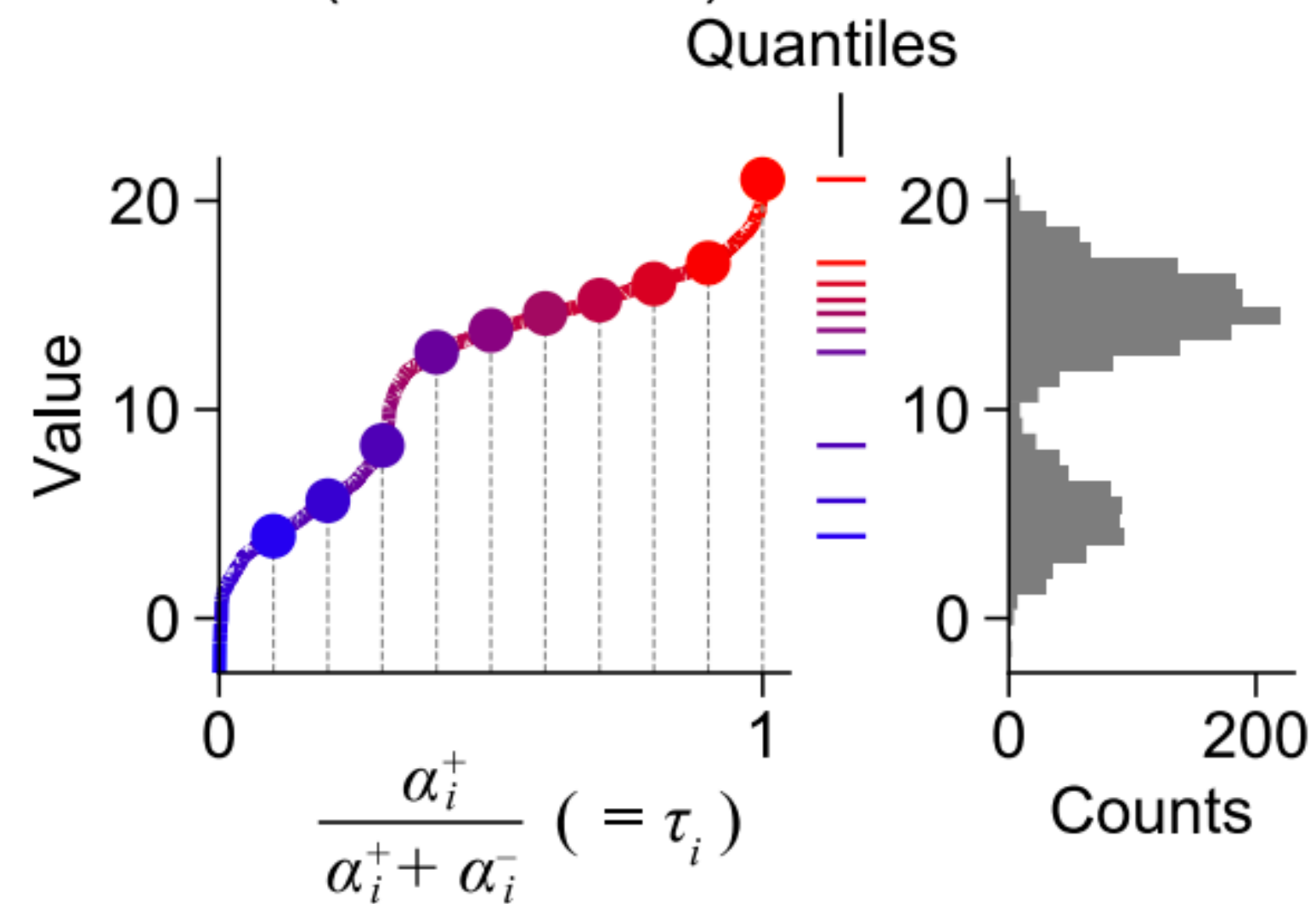
(D)



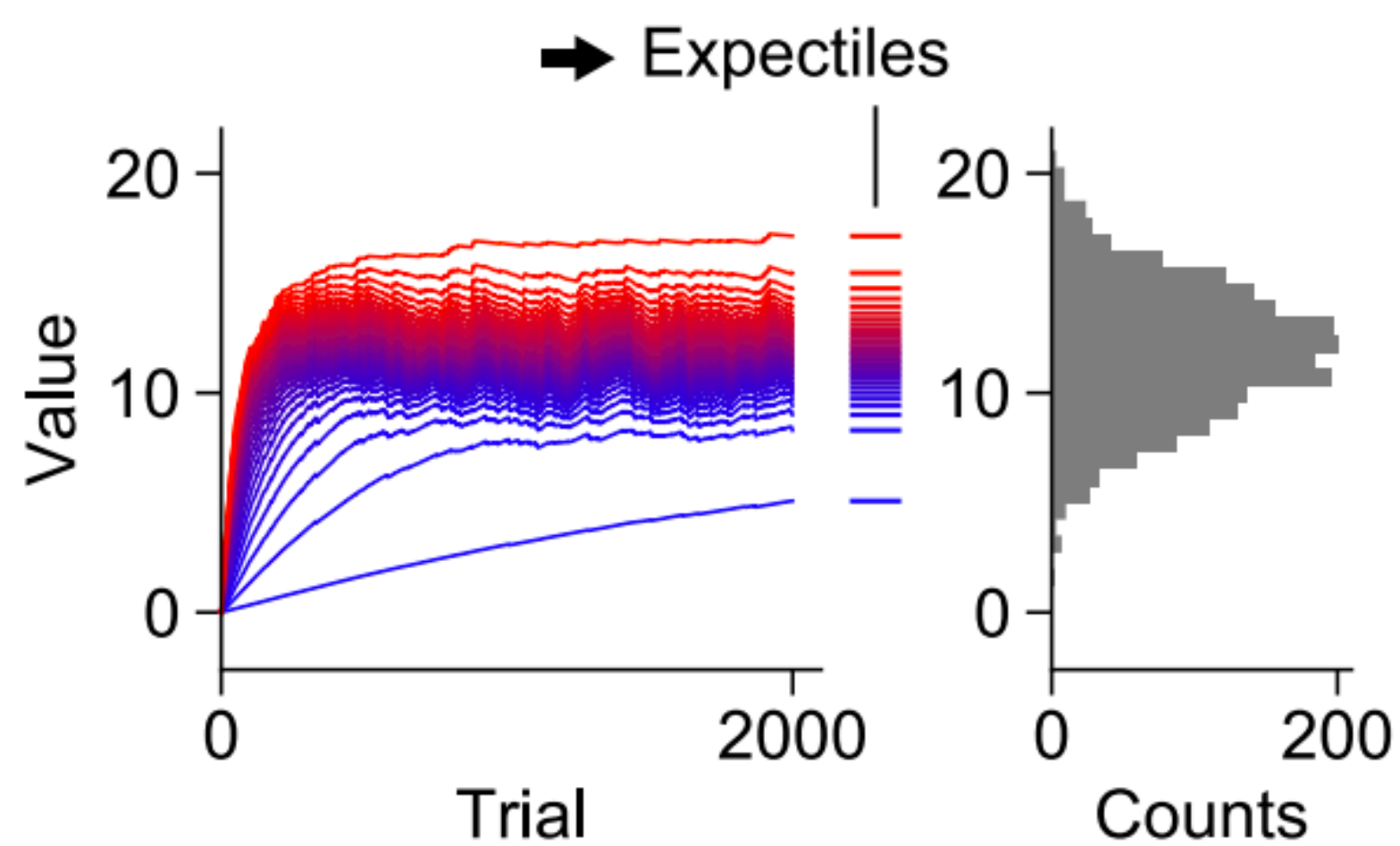
(E) Cumulative distribution function (CDF)



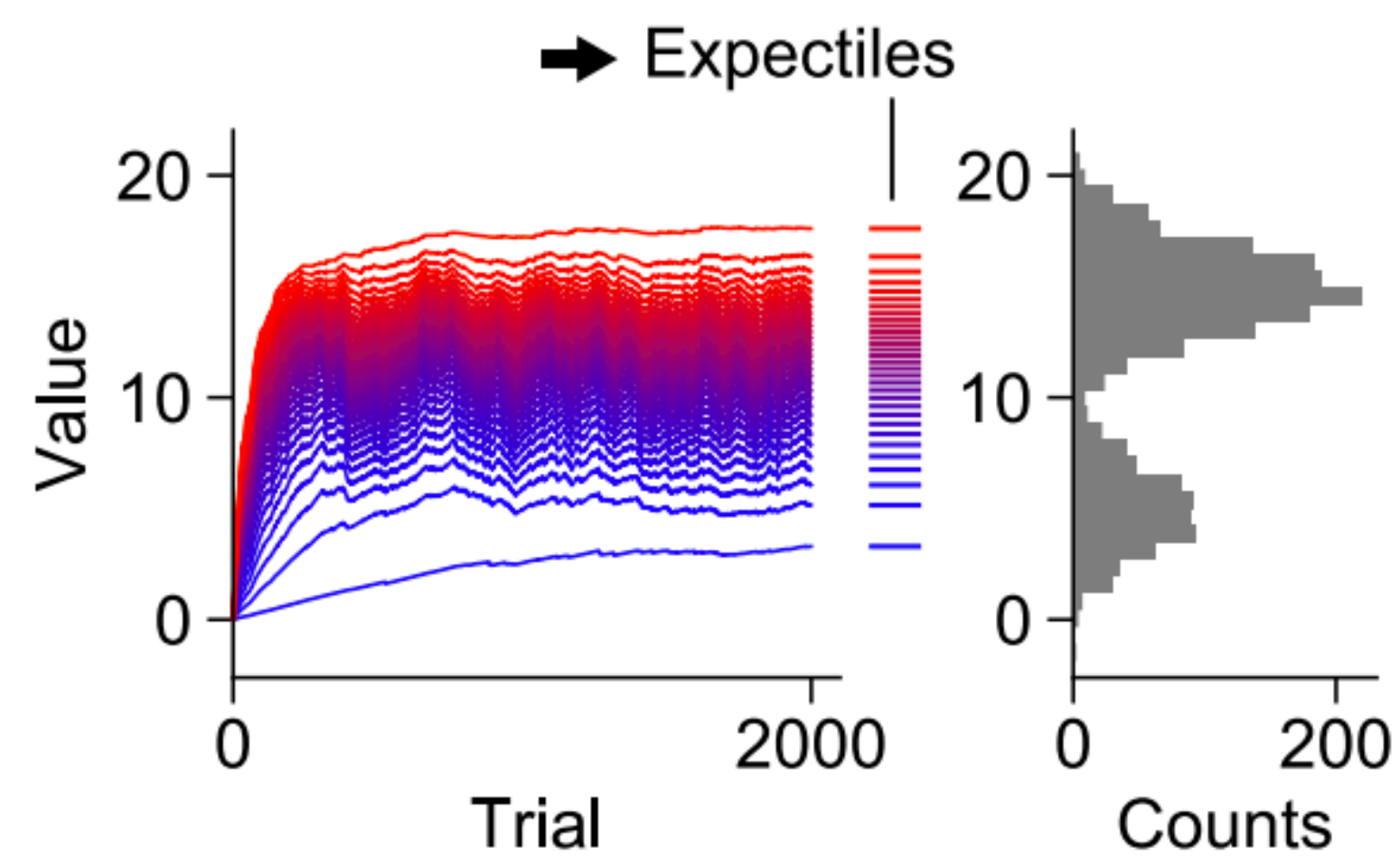
(F) Inverse cumulative distribution function (Inverse CDF)



(G)
$$V_i \leftarrow V_i + \begin{cases} \alpha_i^- \cdot \delta_i & \text{if } \delta_i \leq 0 \\ \alpha_i^+ \cdot \delta_i & \text{if } \delta_i > 0 \end{cases}$$



(H)

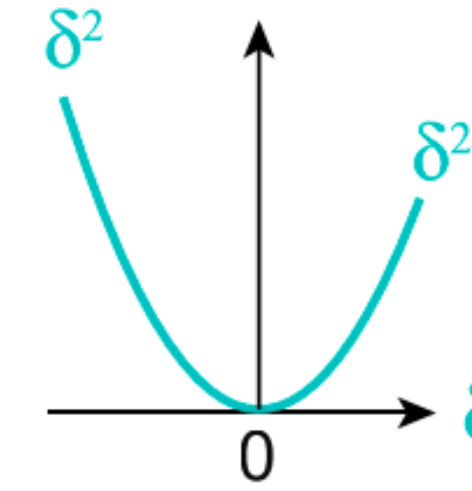
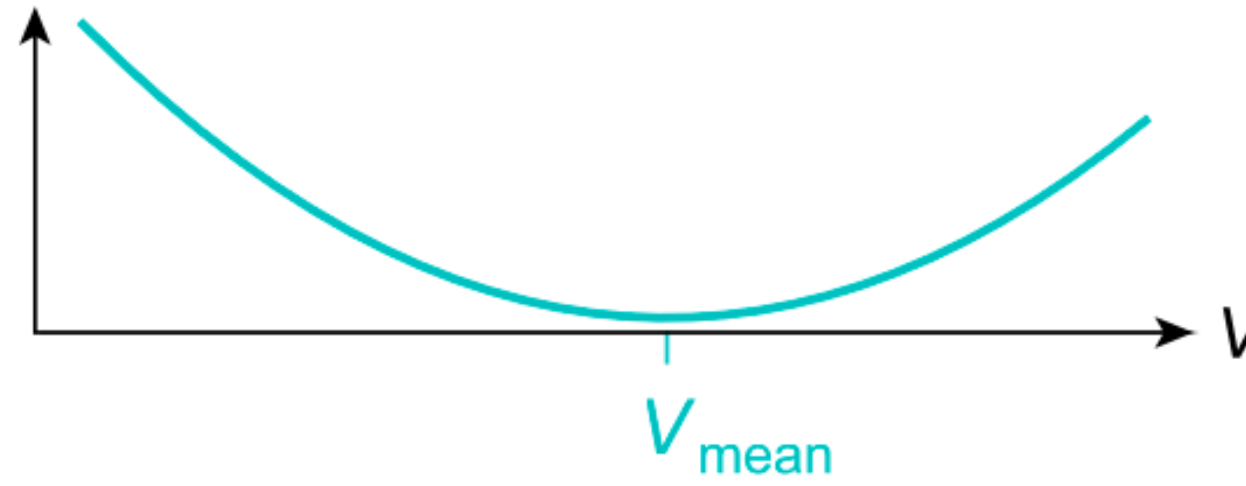


Loss (error) function

$$\delta_n = r_n - V$$

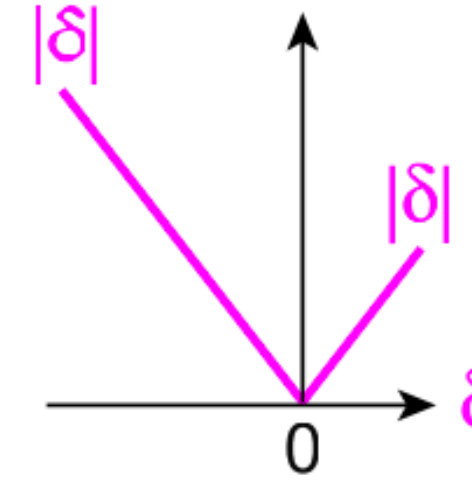
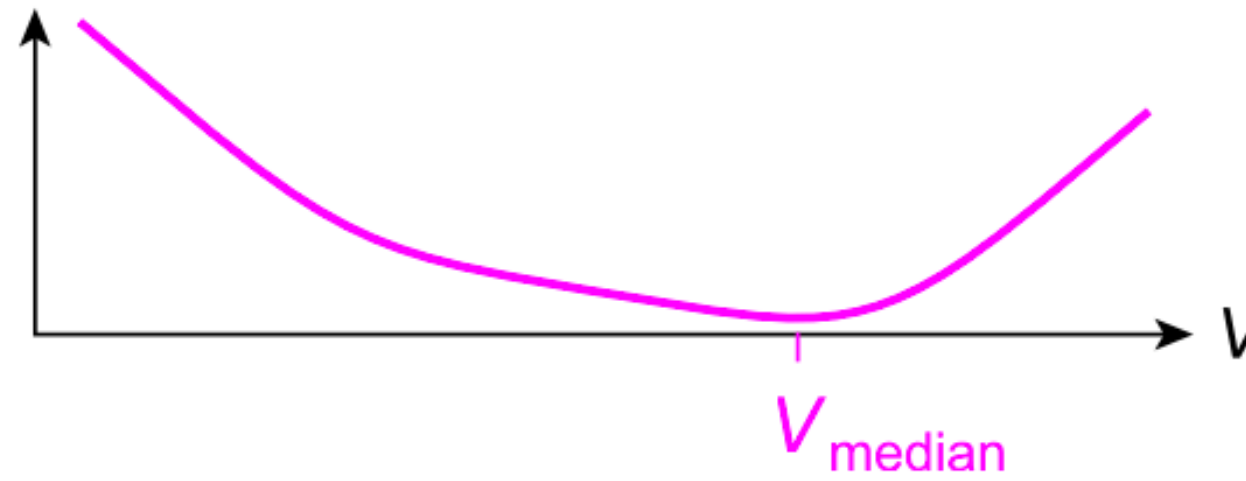
(B) Mean squared error

$$\sum_n \delta_n^2$$



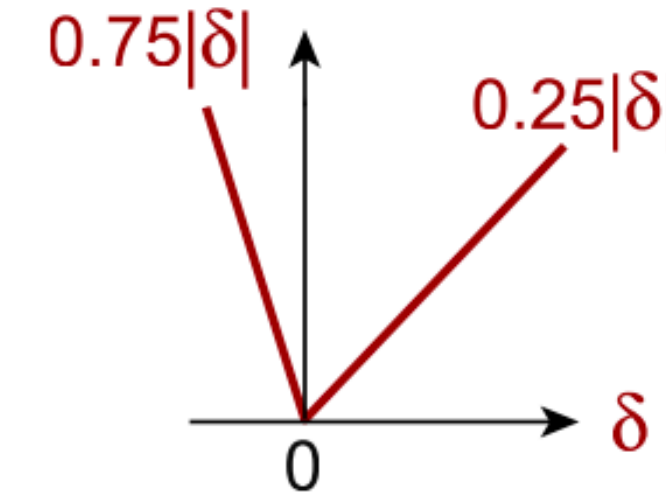
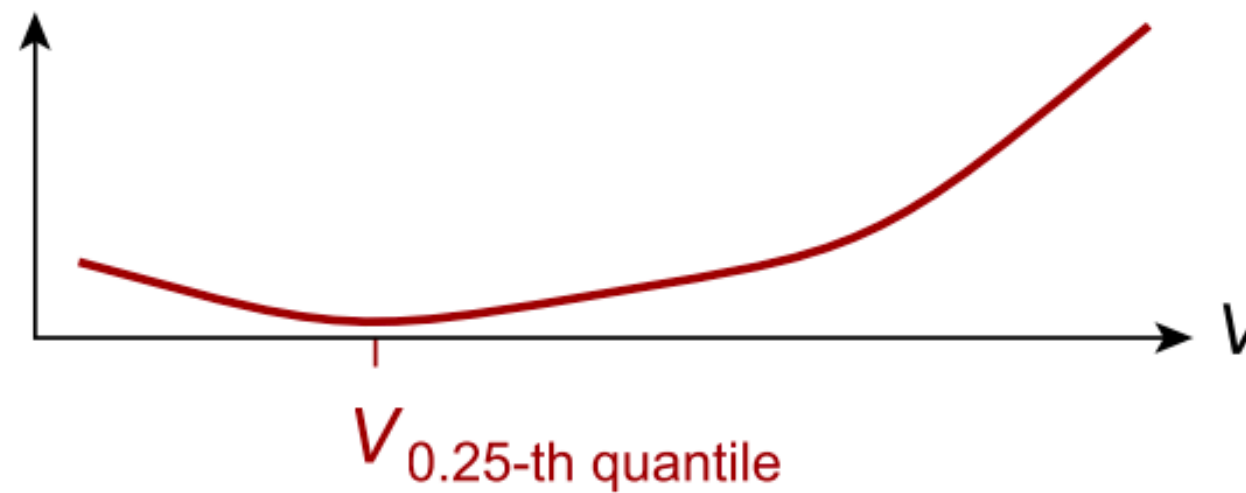
(C) Mean absolute error

$$\sum_n |\delta_n|$$



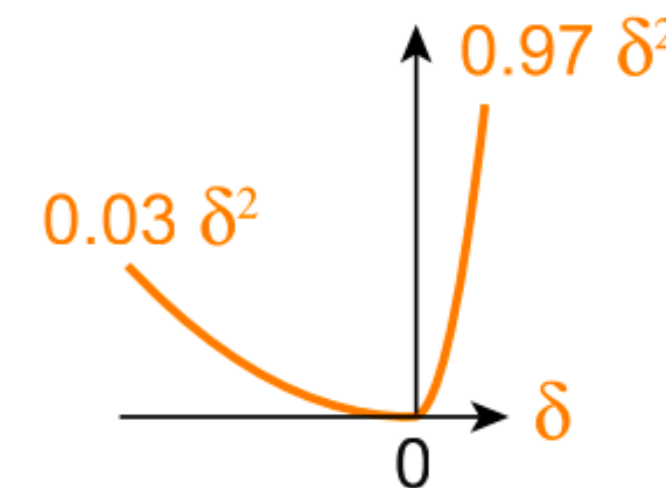
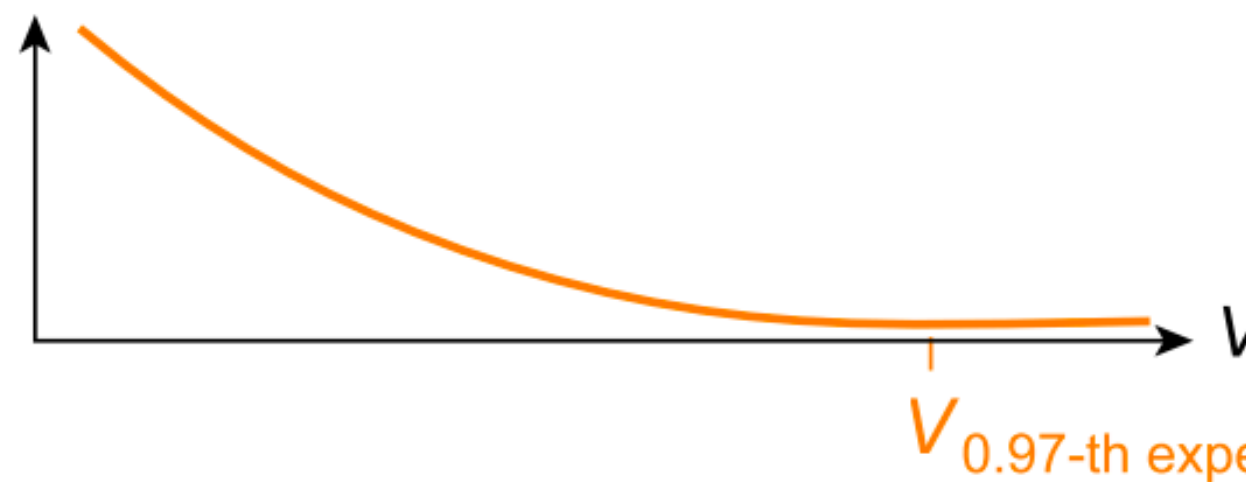
(D) Quantile regression loss

$$\sum_n |\delta_n| \cdot \begin{cases} 0.75 & \text{if } \delta_n \leq 0 \\ 0.25 & \text{if } \delta_n > 0 \end{cases}$$

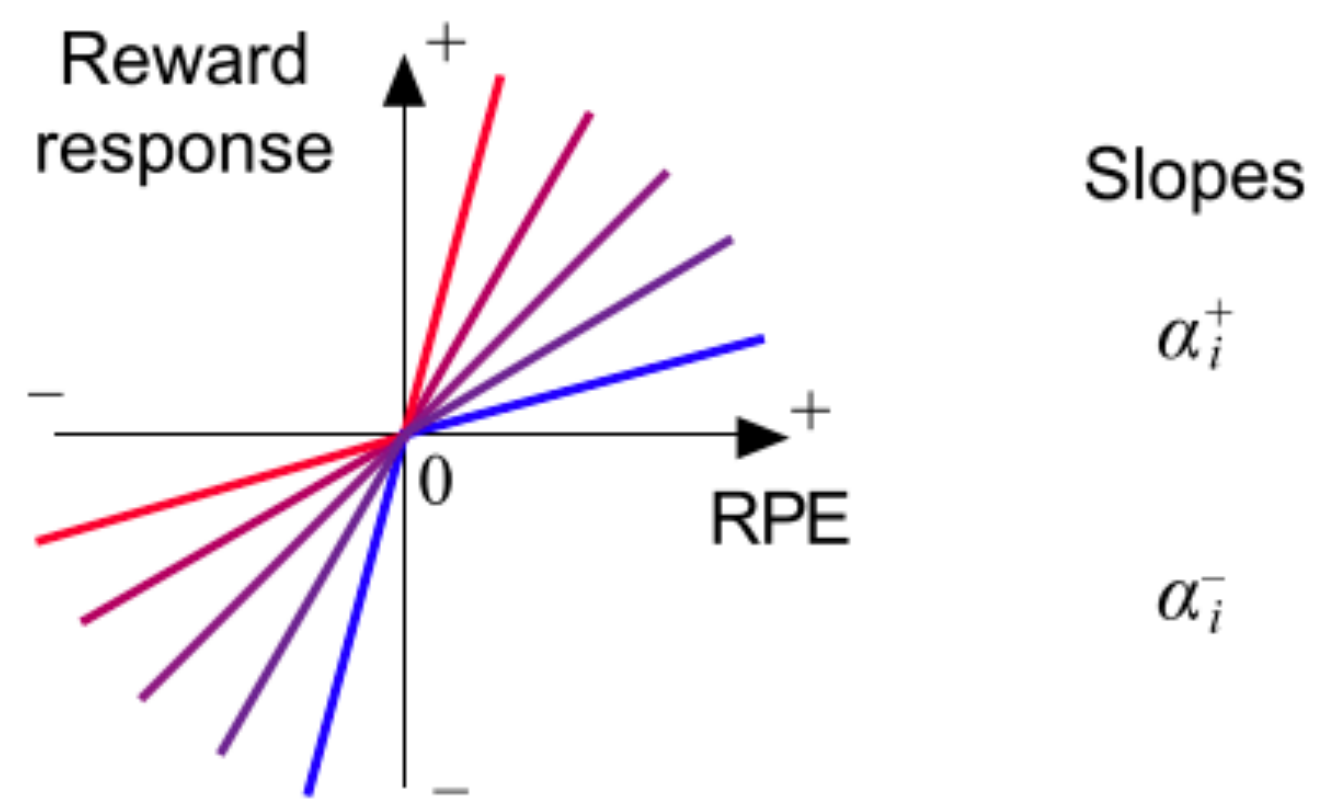


(E) Expectile regression loss

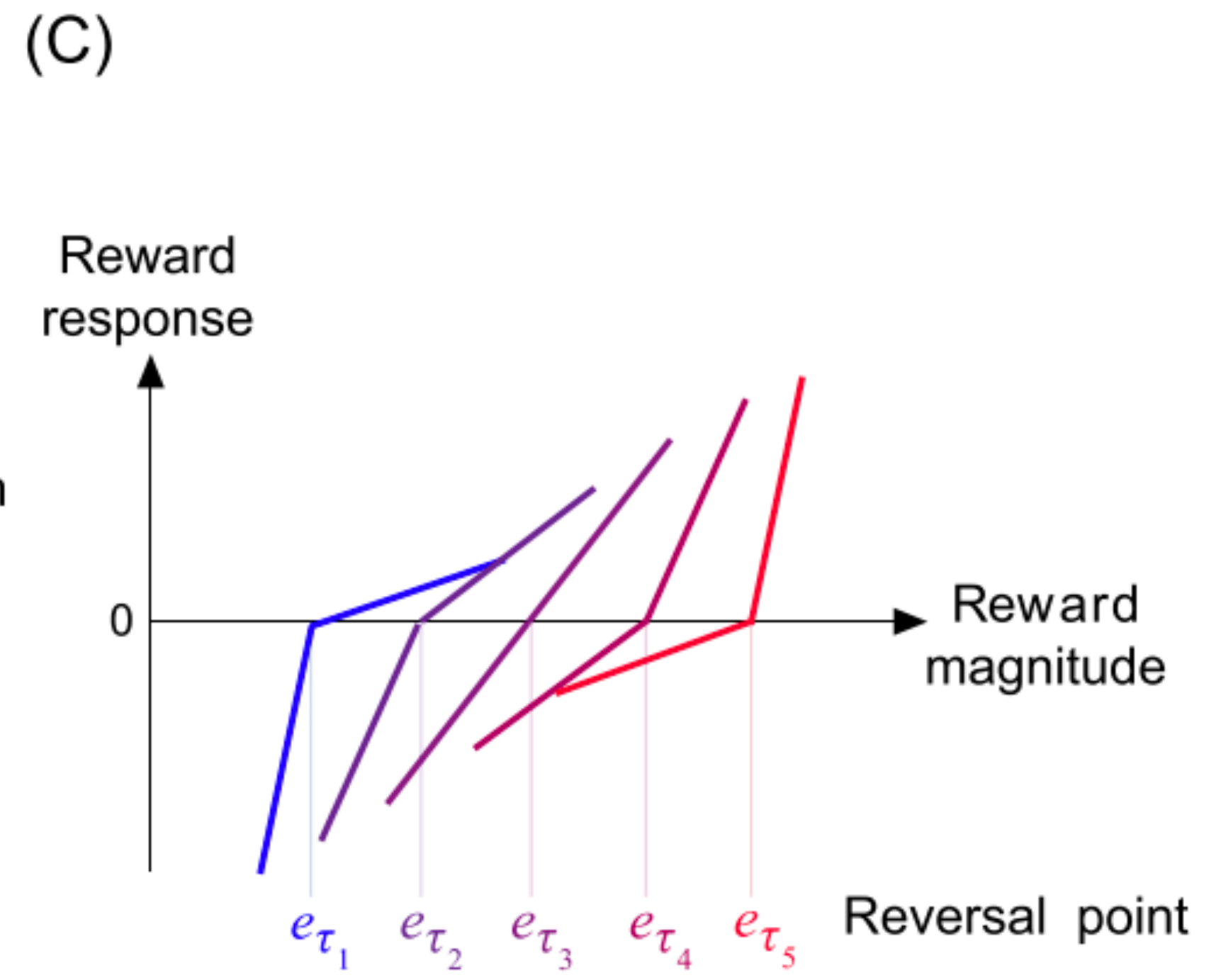
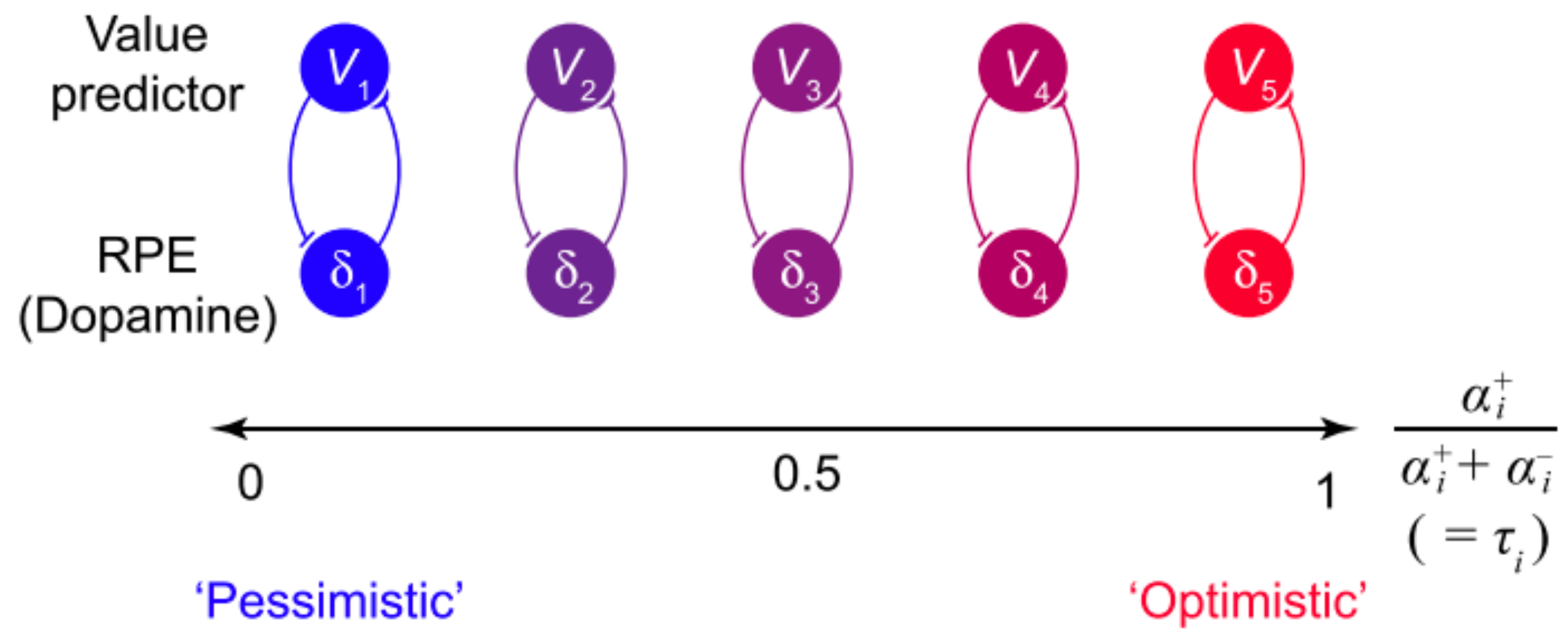
$$\sum_n \delta_n^2 \cdot \begin{cases} 0.03 & \text{if } \delta_n \leq 0 \\ 0.97 & \text{if } \delta_n > 0 \end{cases}$$



(A) (1) Variability in RPE encoding

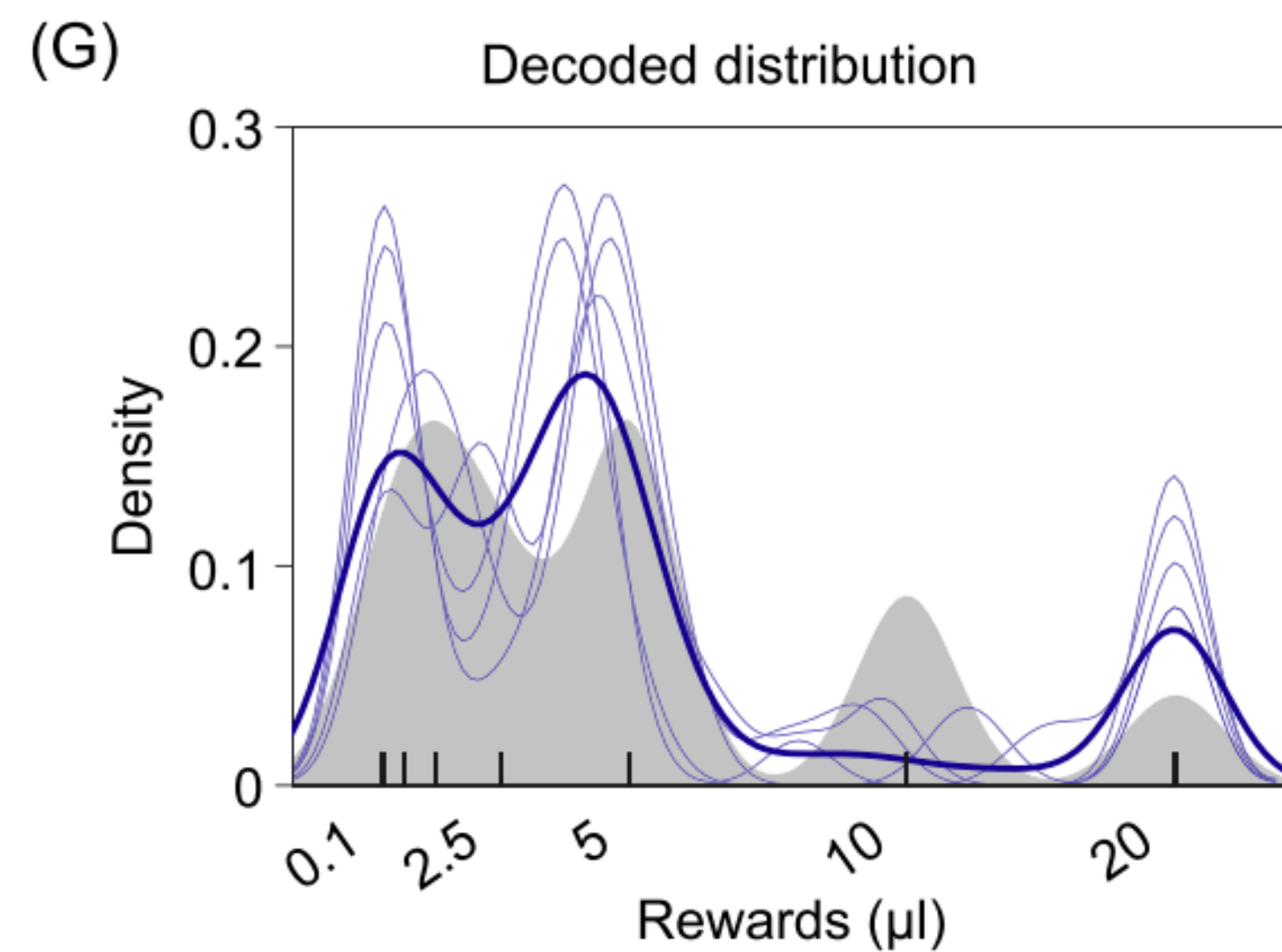
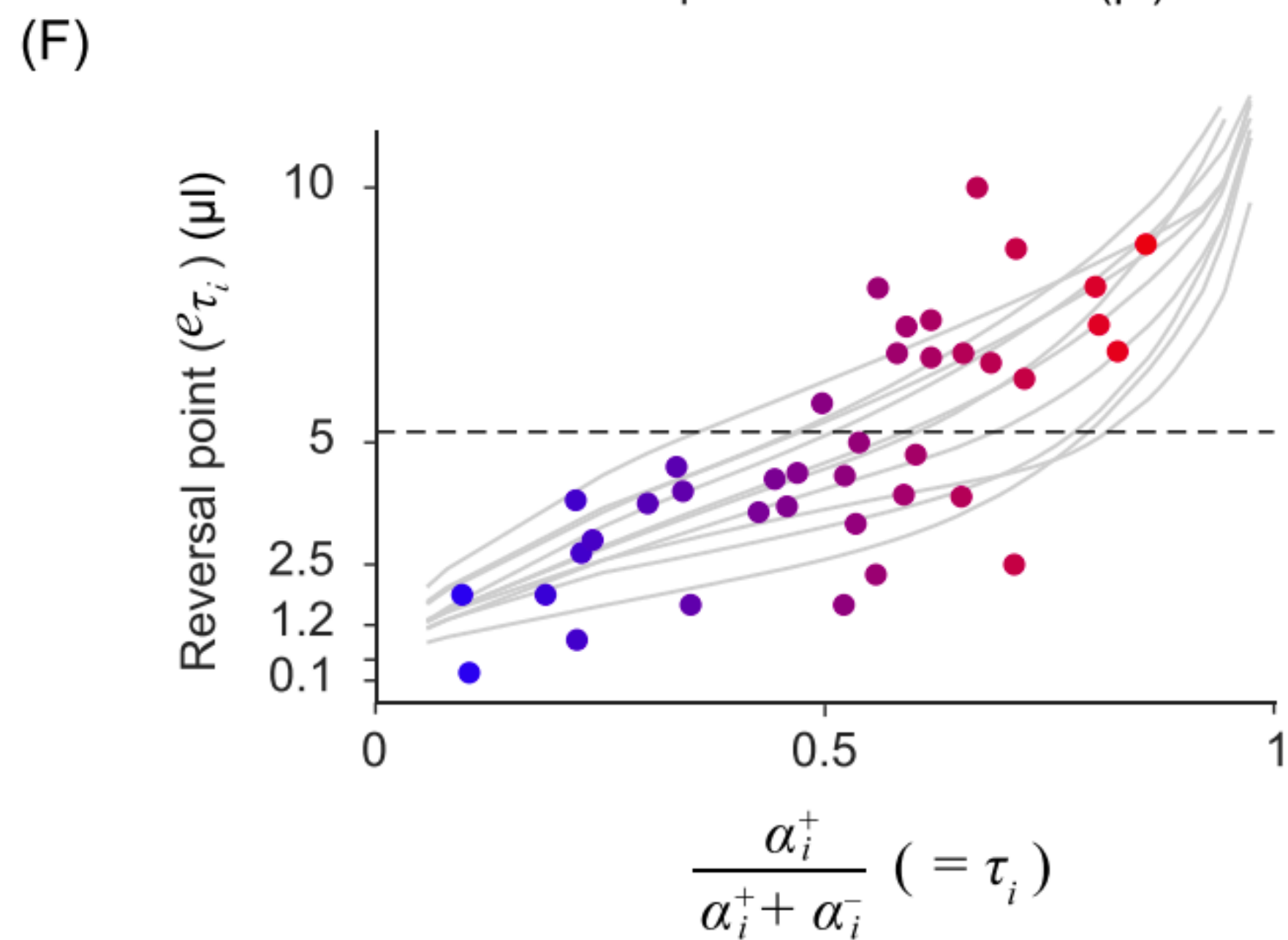
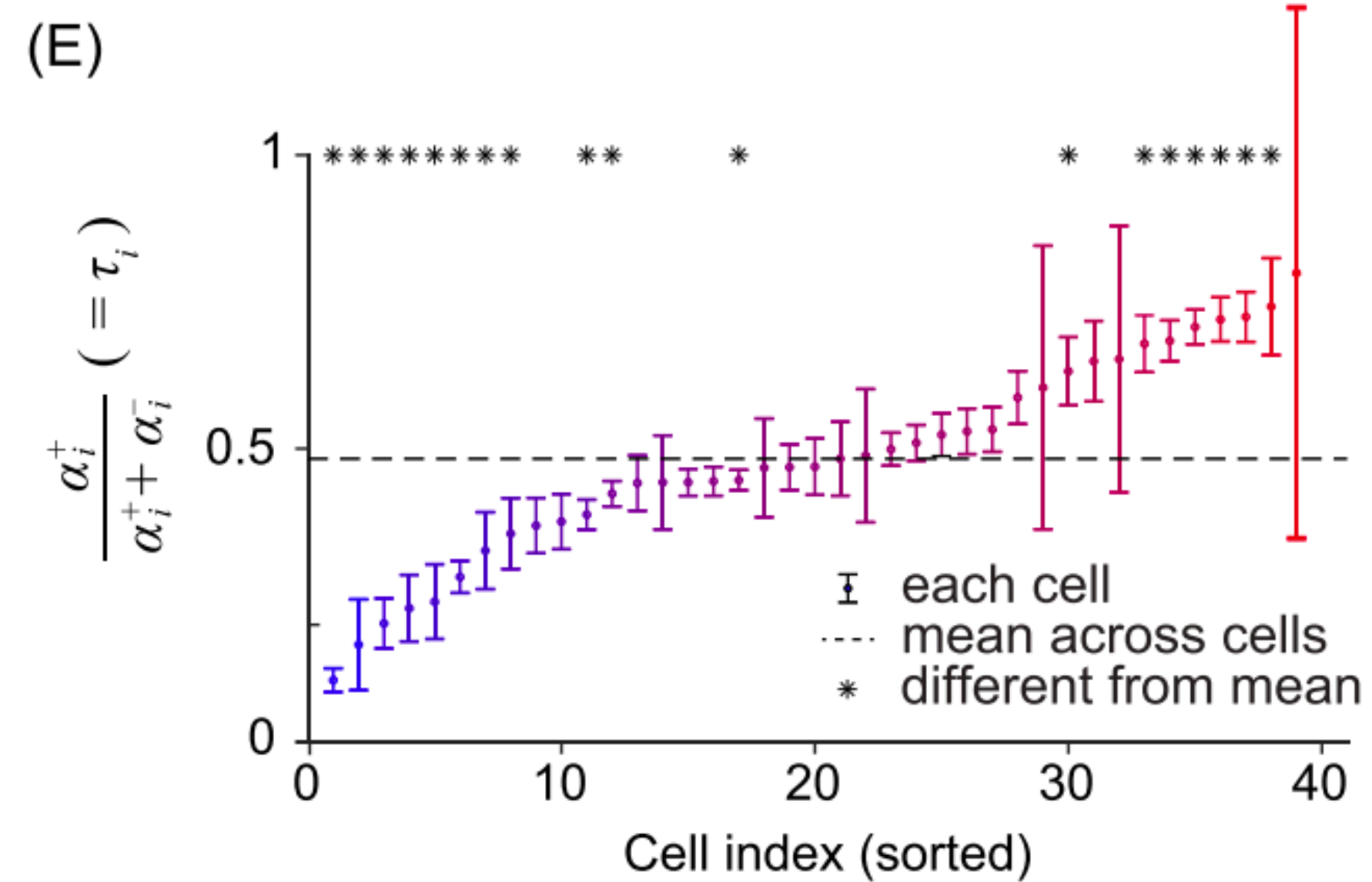
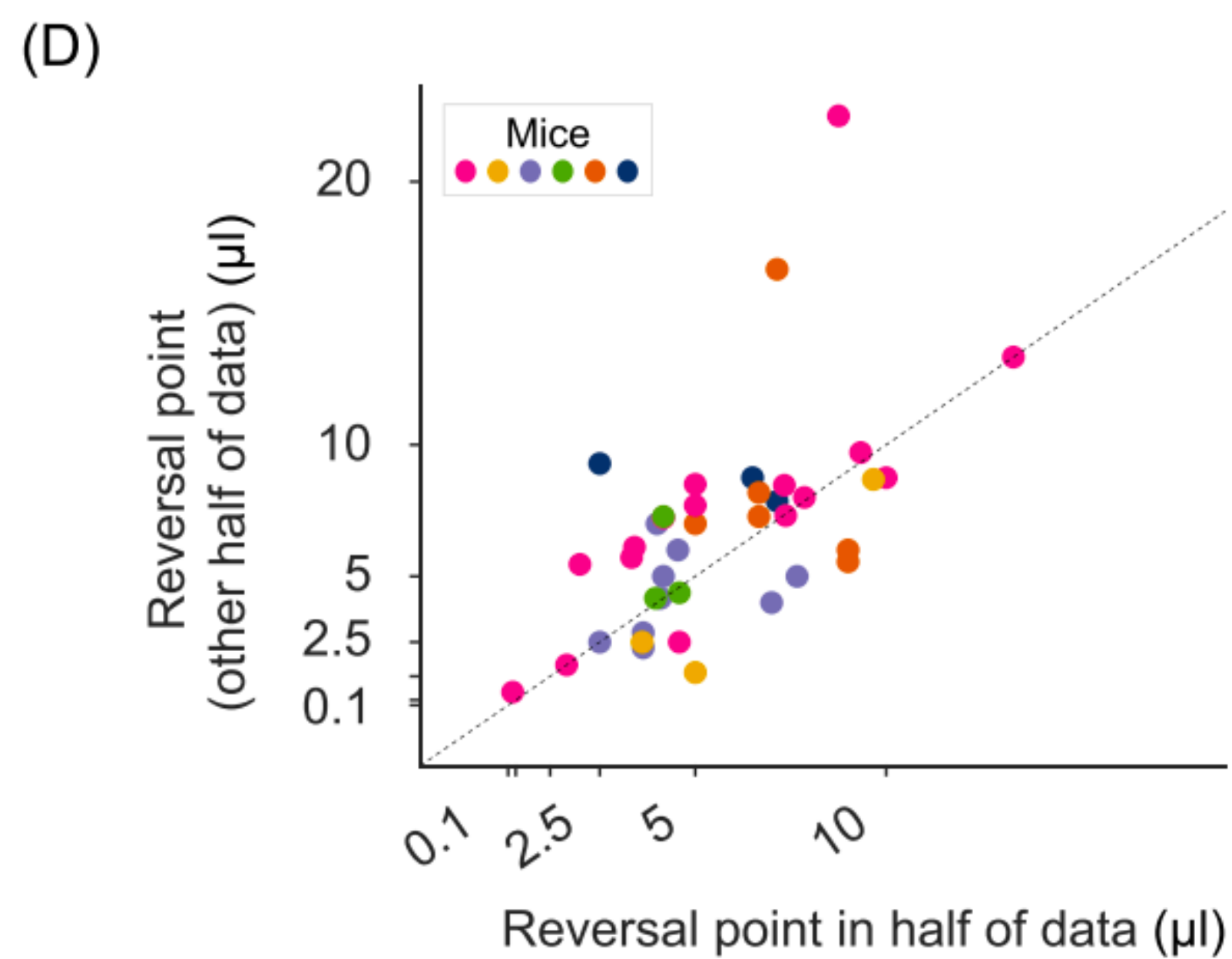


(B) (2) Family of value predictors/RPEs



· Reversal point \approx τ -th expectile or quantile

$$\tau_i = \frac{\alpha_i^+}{\alpha_i^+ + \alpha_i^-}$$



Distributional TD learning

- In distributional TD learning, one needs to sample the value function from the distribution
- It is easy in the quantile code by taking linear combination
- Currently, the performance boost from distributional RL comes from its benefit on learning the representation. How does knowing a distribution involves in decision making is another question.