

Memory and the statistical structure of the world

Samuel Gershman

Department of Brain and
Cognitive Sciences, MIT

July 2014

Marco Polo



Marco Polo



Is this a funny-looking unicorn?



Big picture

- When do we modify old memories, and when do we create new ones?

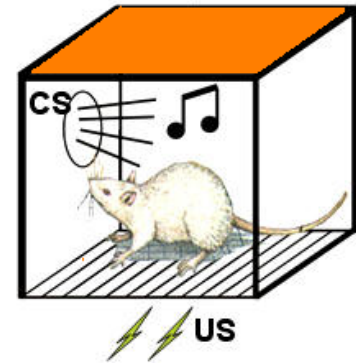
Big picture

- When do we modify old memories, and when do we create new ones?
- This question can be answered within a probabilistic computational framework:
we create new memories when we infer new latent causes in our environment

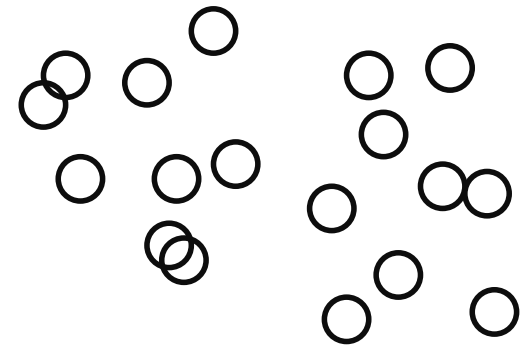
Big picture

- When do we modify old memories, and when do we create new ones?
- This question can be answered within a probabilistic computational framework:
we create new memories when we infer new latent causes in our environment
- This principle has deep explanatory power across multiple domains

- Classical conditioning



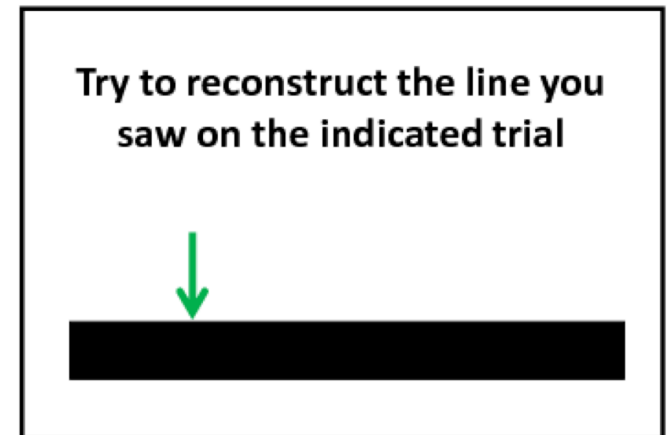
- Perceptual estimation



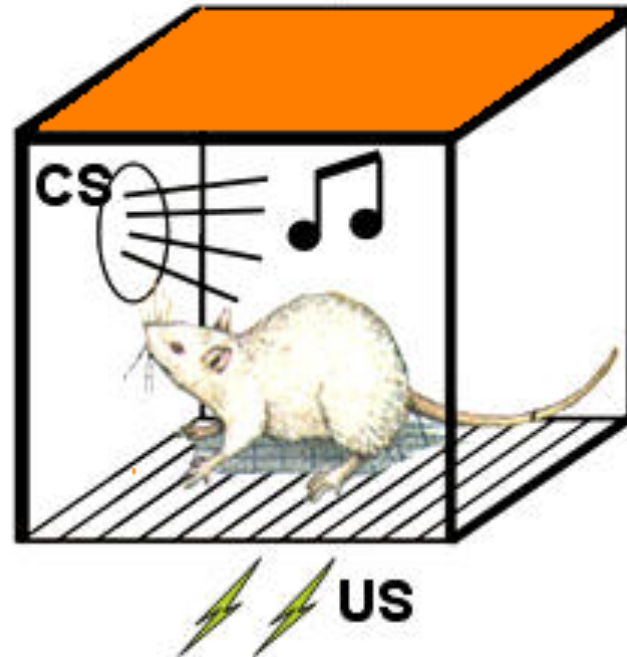
How many circles?

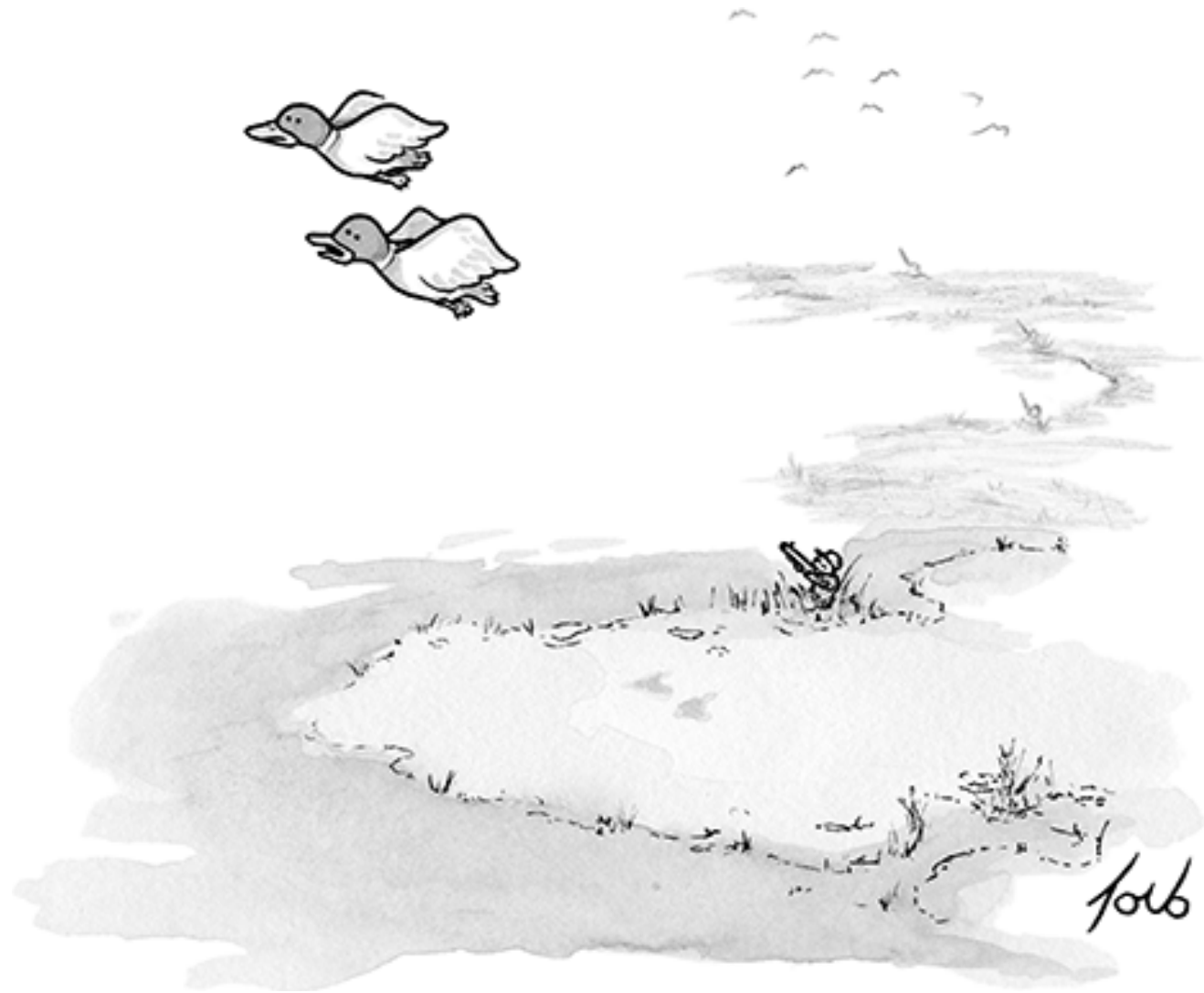
Reconstruction trial

- Reconstructive memory



What do animals learn during classical conditioning?

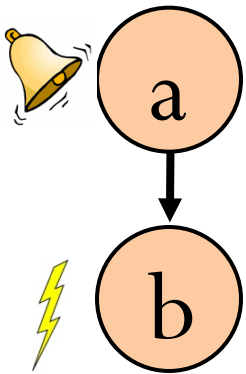




“It’s that time of year when guys randomly explode.”

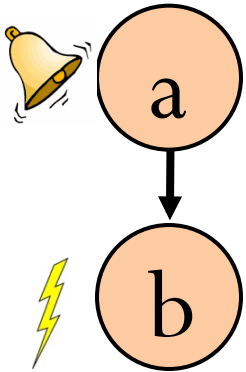
Some possibilities

Tone (a) causes
shock (b)

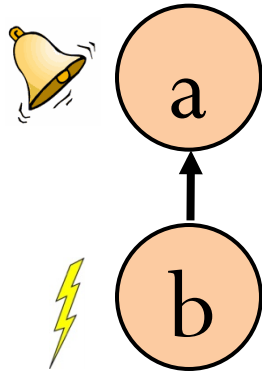


Some possibilities

Tone (a) causes
shock (b)

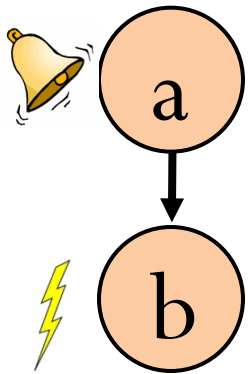


Shock causes
tone

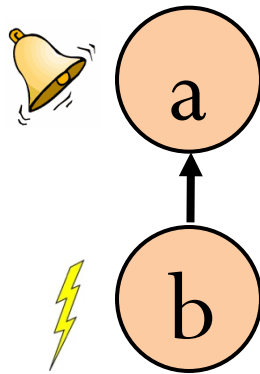


Some possibilities

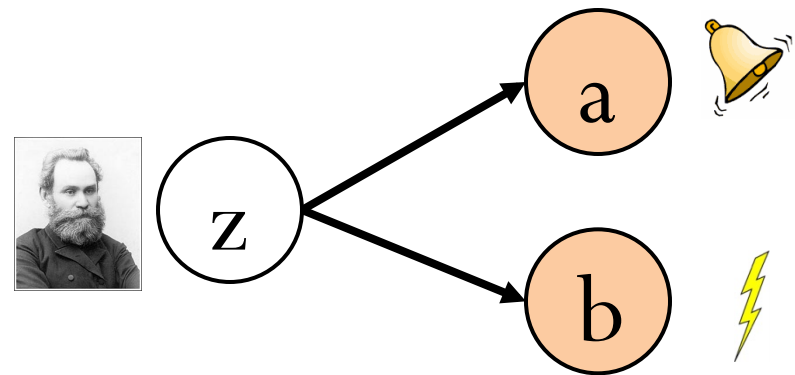
Tone (a) causes
shock (b)



Shock causes
tone

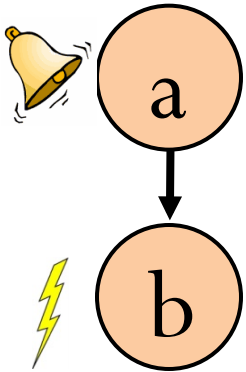


Something else (z) causes
both tone and shock

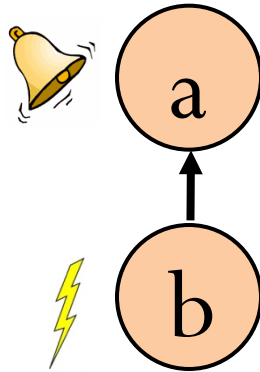


Some possibilities

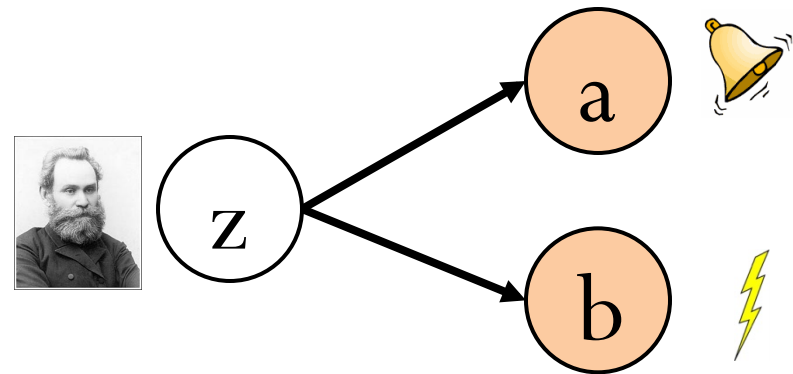
Tone (a) causes
shock (b)



Shock causes
tone



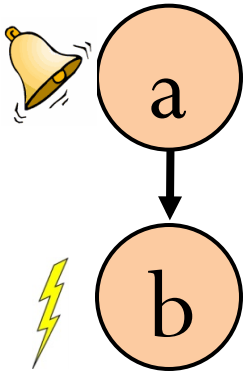
Something else (z) causes
both tone and shock



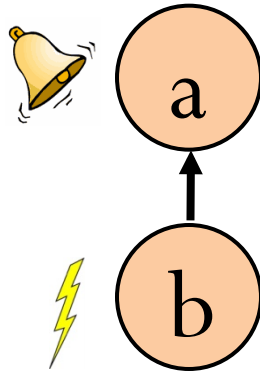
Too constrained

Some possibilities

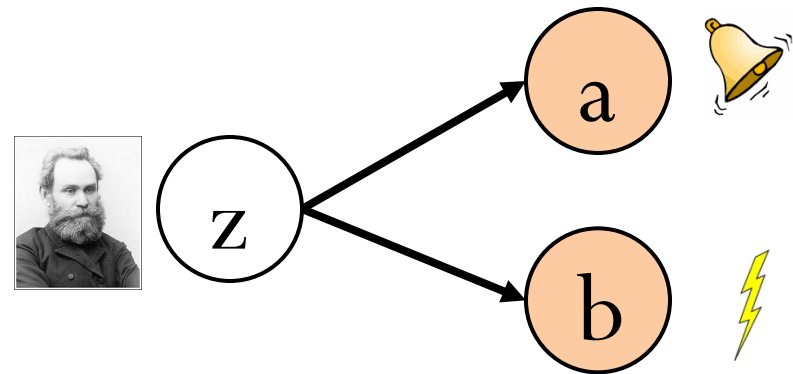
Tone (a) causes
shock (b)



Shock causes
tone



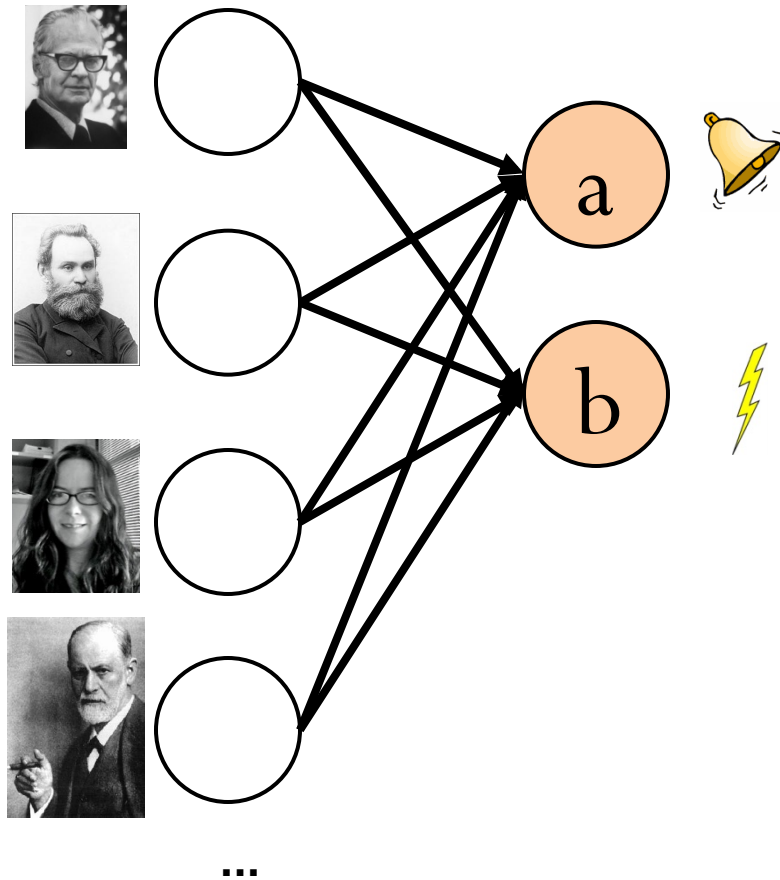
Something else (z) causes
both tone and shock



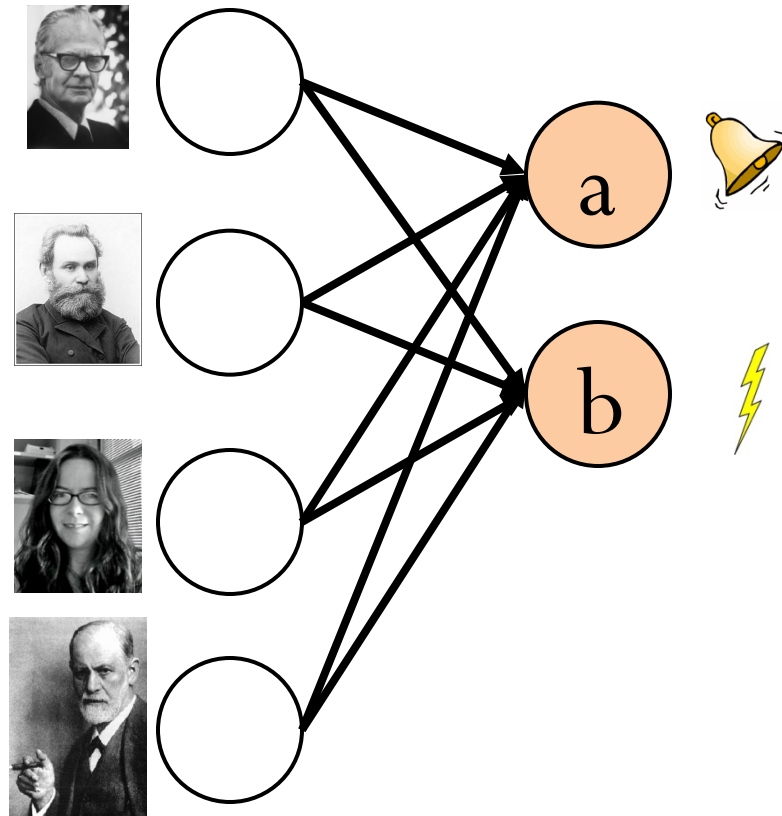
Too constrained

Too flexible?

Too flexible?



Too flexible?

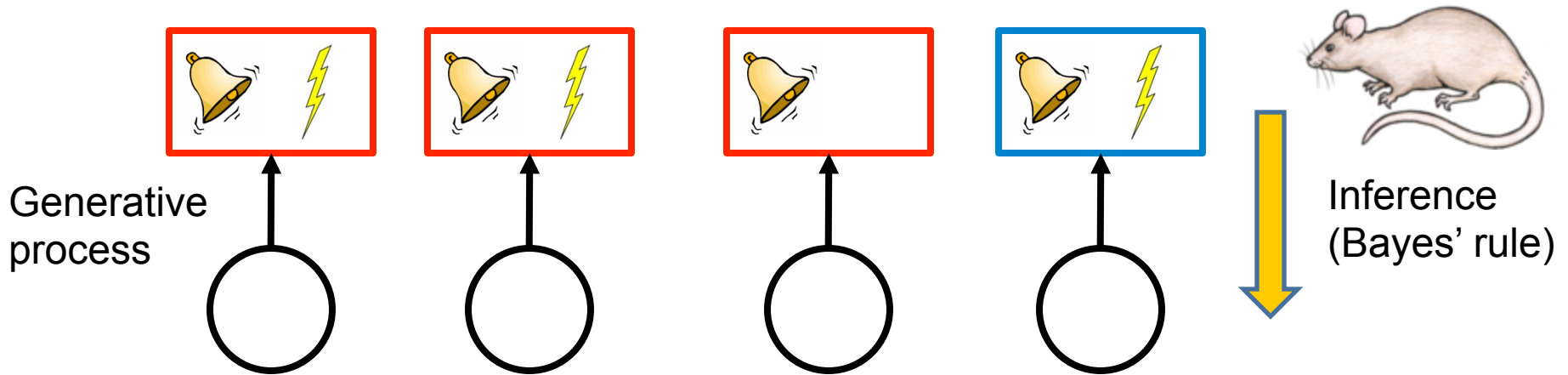


Hypothesis: Animals assume a generative model in which (1) the number of latent causes is unbounded, and (2) a small number of latent causes is more likely *a priori*.

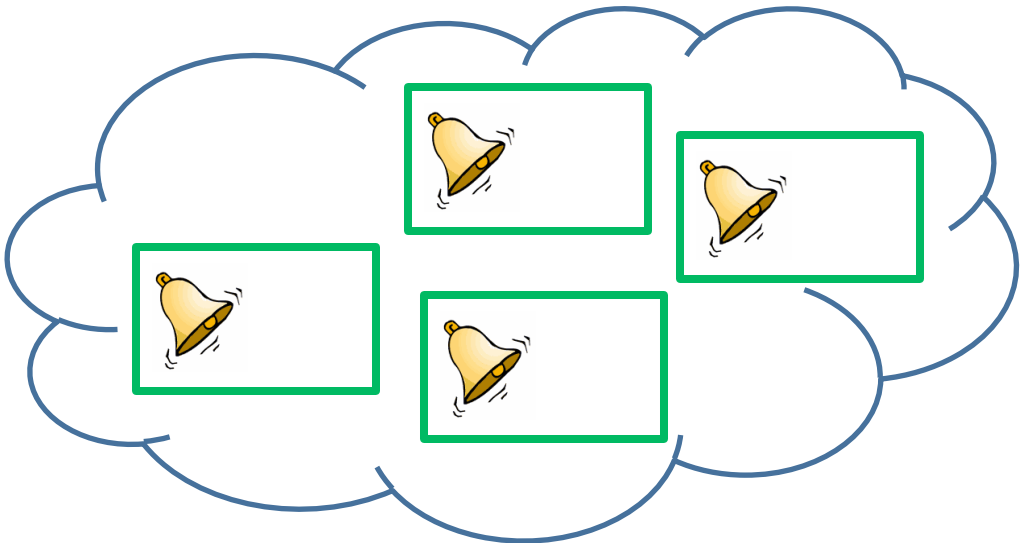
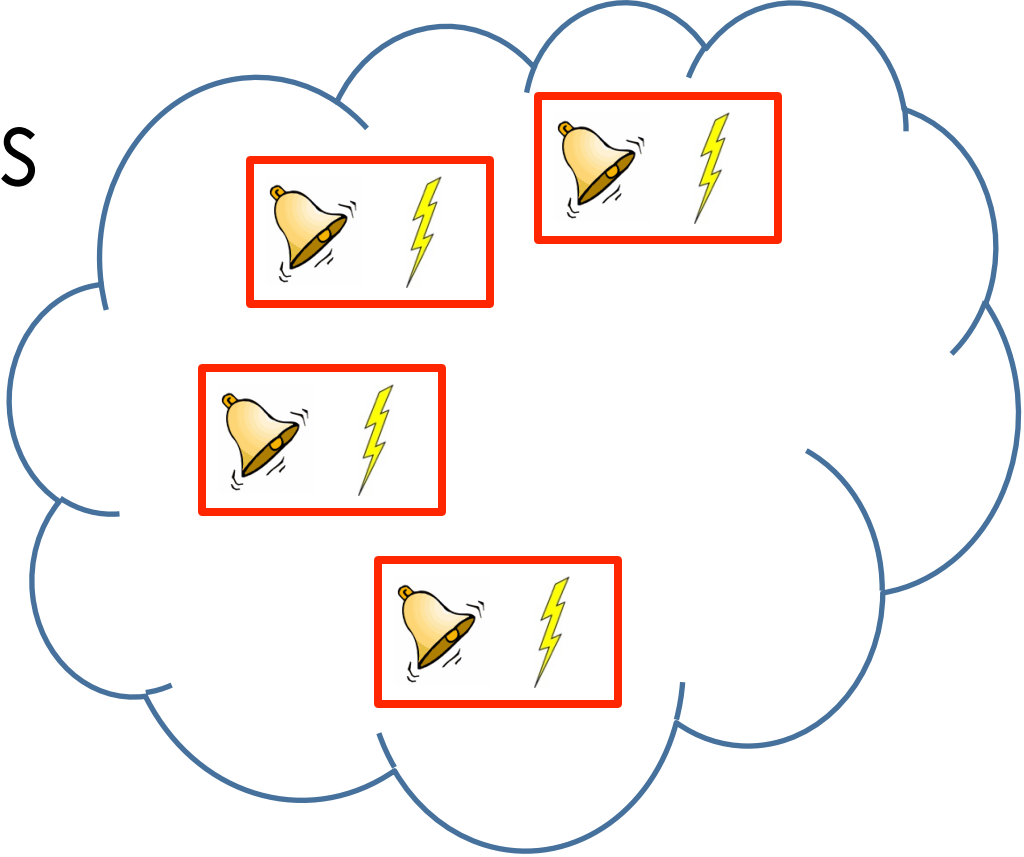
Inverting the generative model

Bayes' rule inverts generative model to infer latent causes:

$$P(\text{cause} \mid \text{data}) \propto P(\text{data} \mid \text{cause})P(\text{cause})$$

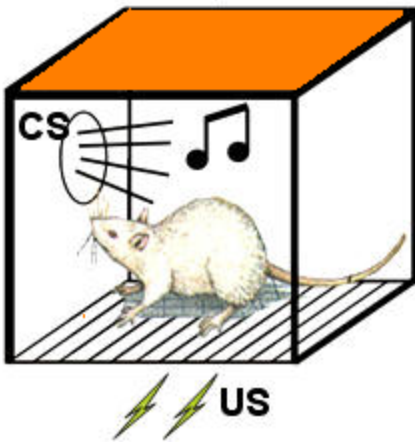


Conditioning as clustering



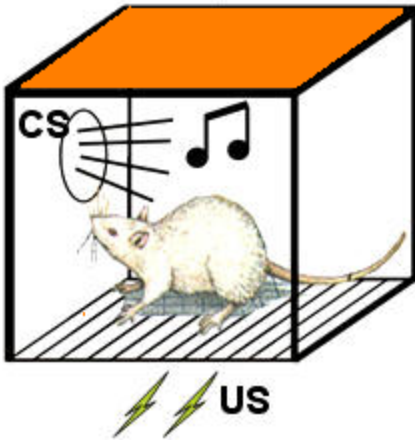
Case study: renewal

Acquisition (box A)

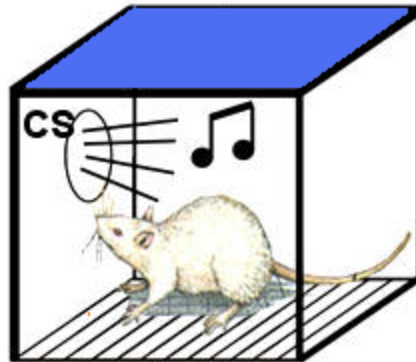


Case study: renewal

Acquisition (box A)

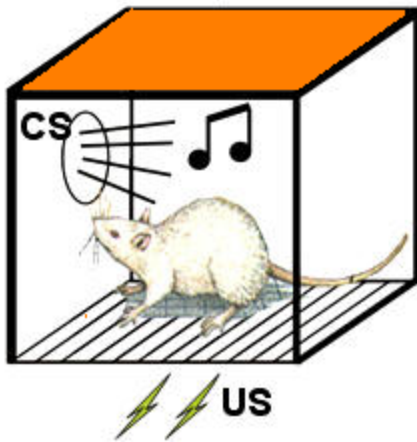


Extinction (box B)

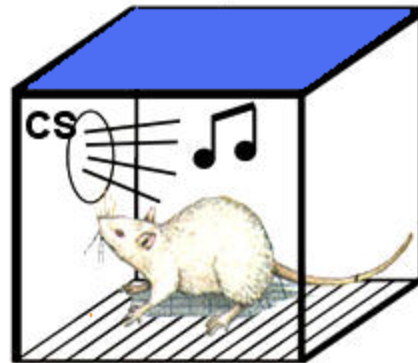


Case study: renewal

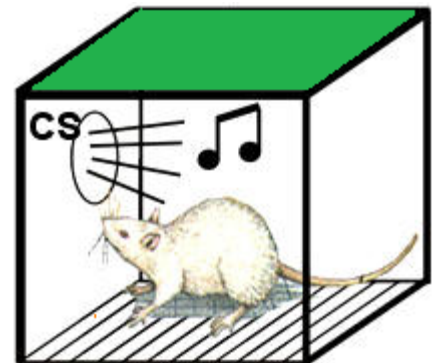
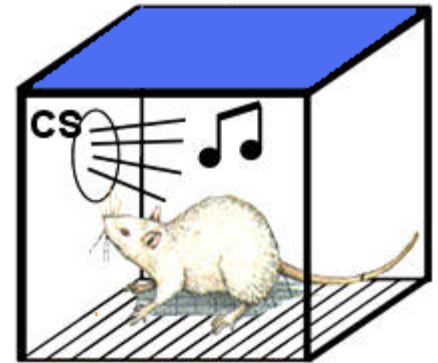
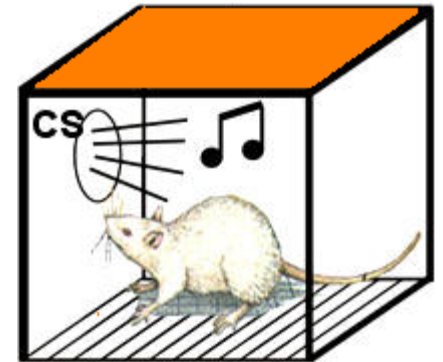
Acquisition (box A)



Extinction (box B)

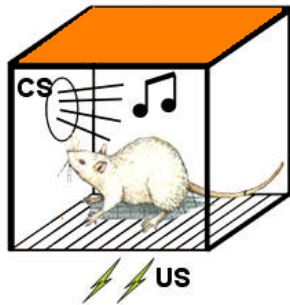


Test

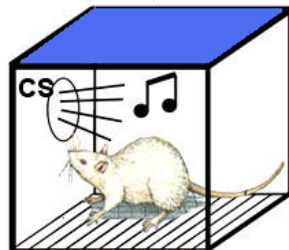


Conditioned responding is renewed!

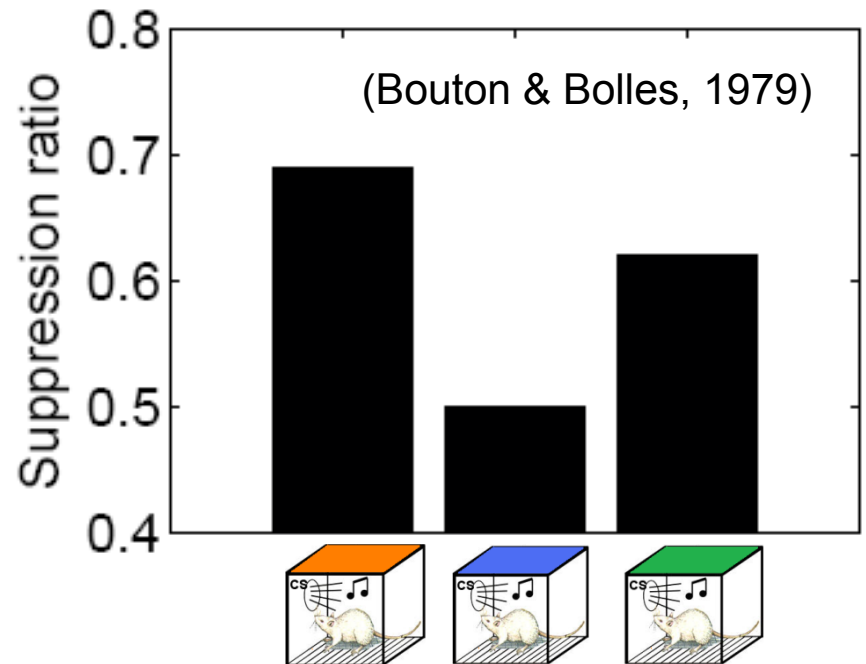
Acquisition



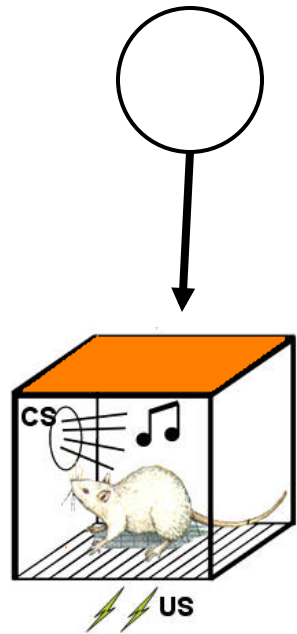
Extinction



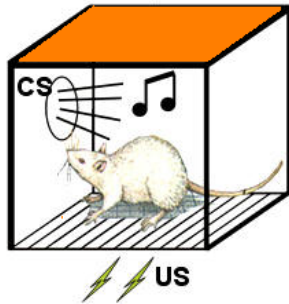
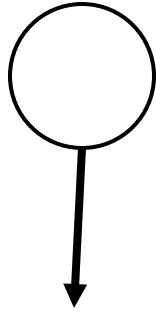
Experimental data



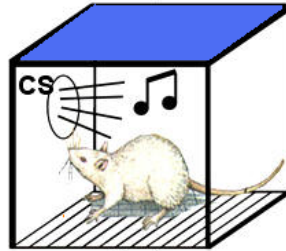
The rat hasn't unlearned its conditioned response; it has **learned something new**.



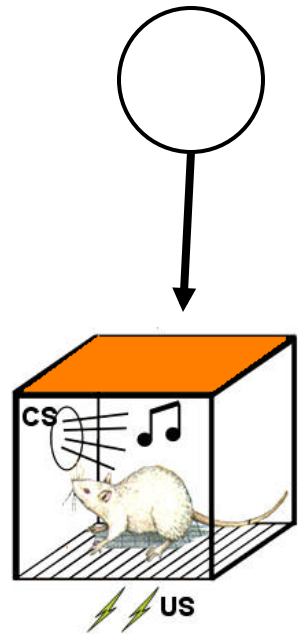
Acquisition



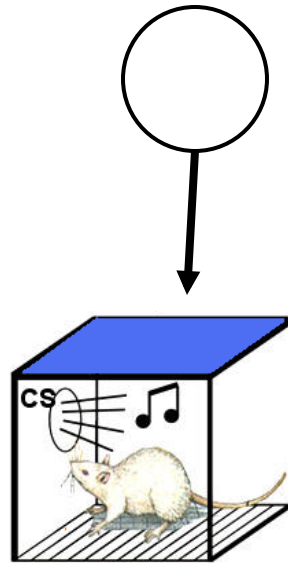
Acquisition



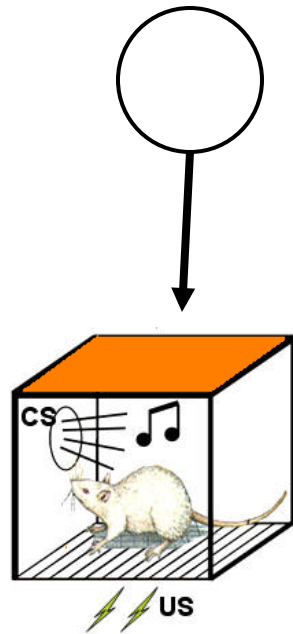
Extinction



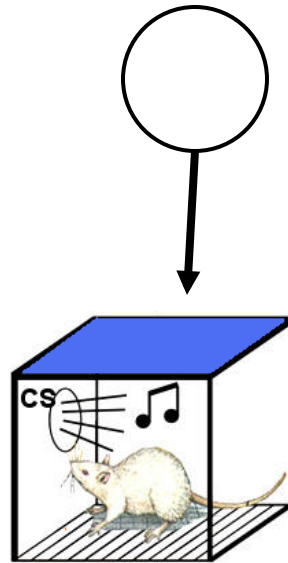
Acquisition



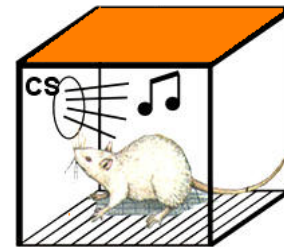
Extinction



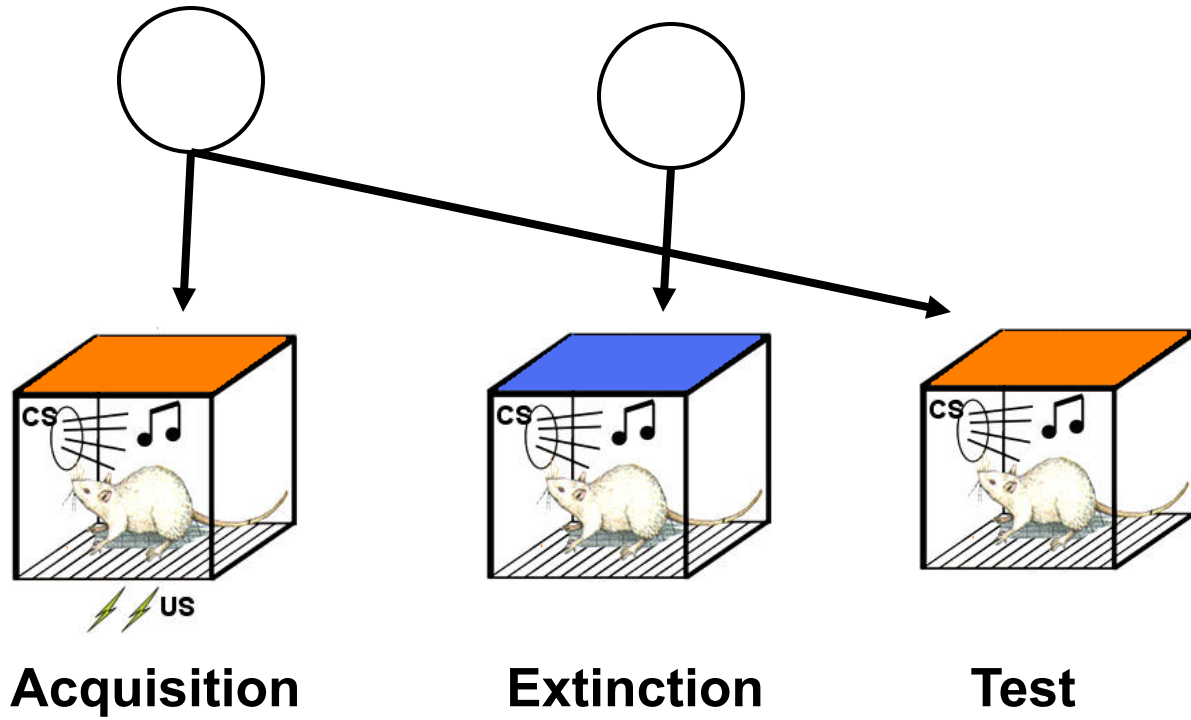
Acquisition



Extinction



Test



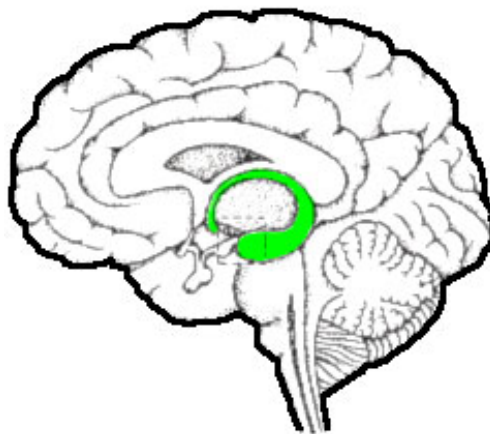
Acquisition

Extinction

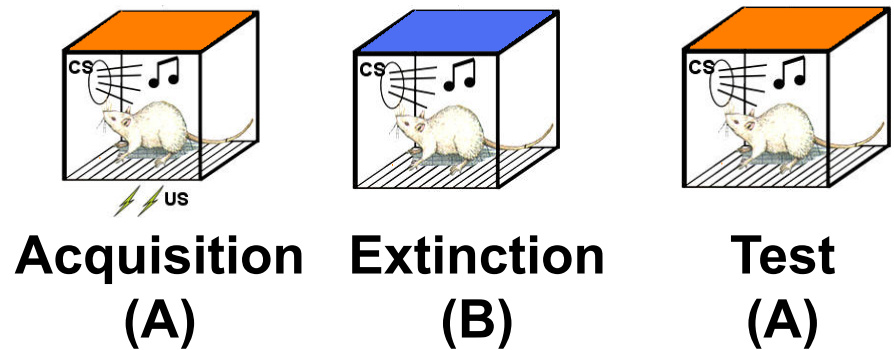
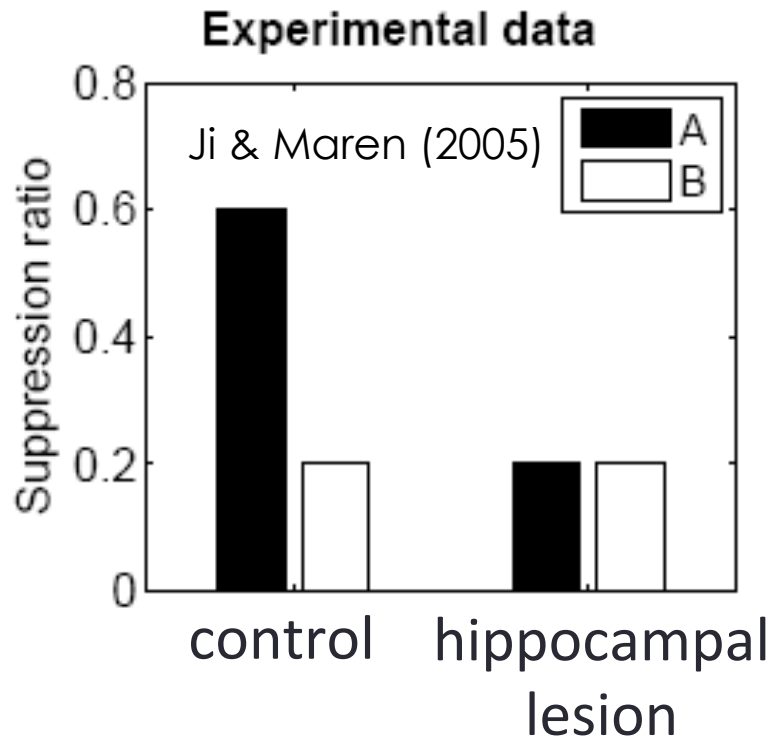
Test

Clustering in the brain

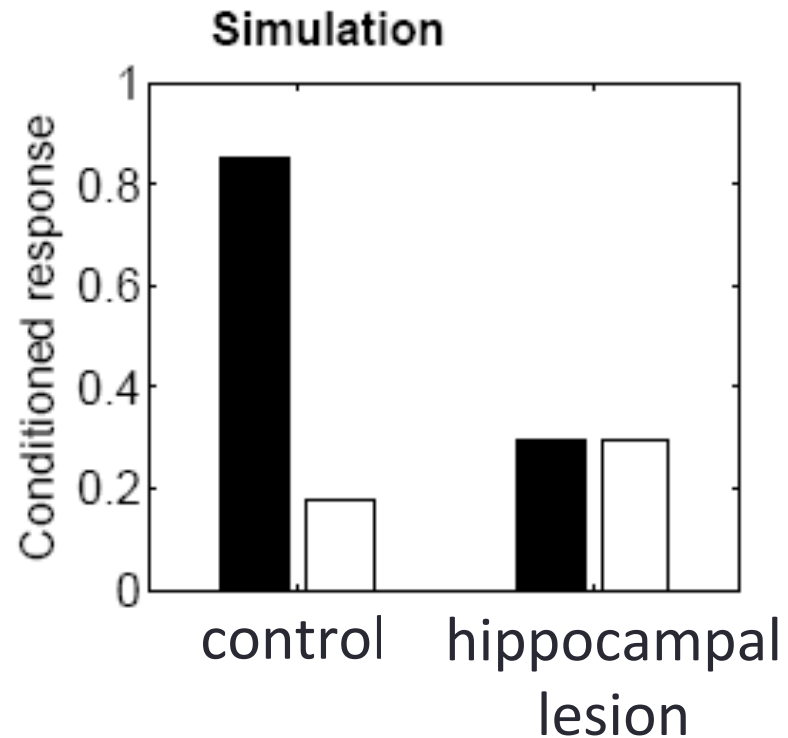
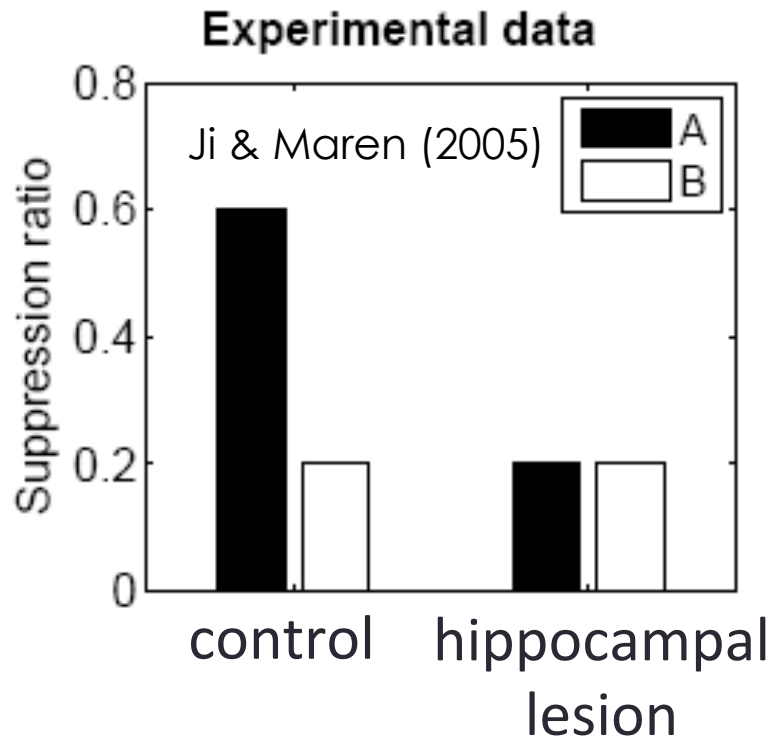
Hippocampus supports the ability to flexibly infer new latent causes



Pre-training lesions of hippocampus abolish renewal

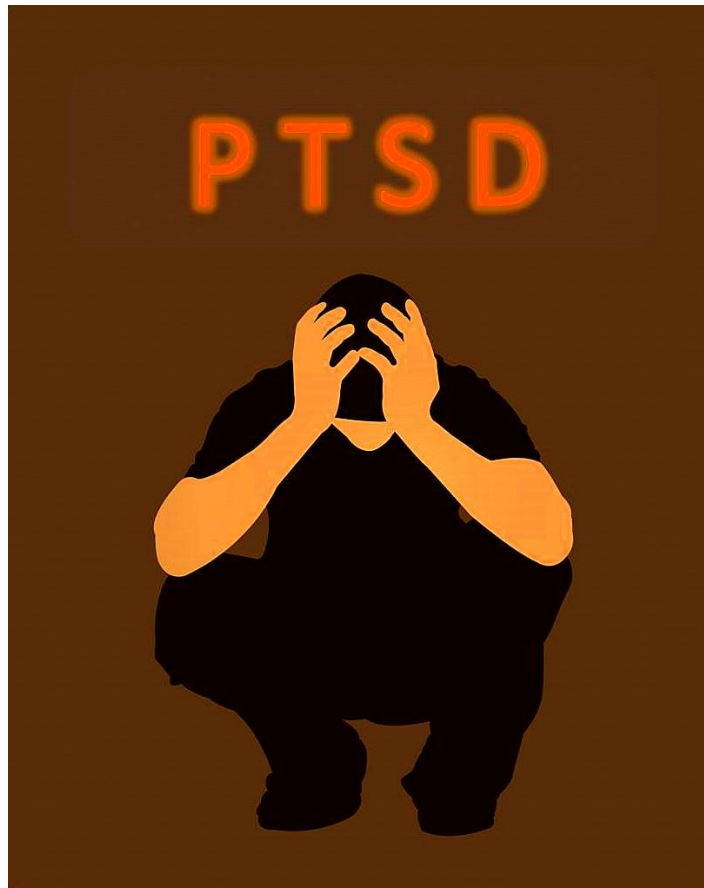


Pre-training lesions of hippocampus abolish renewal

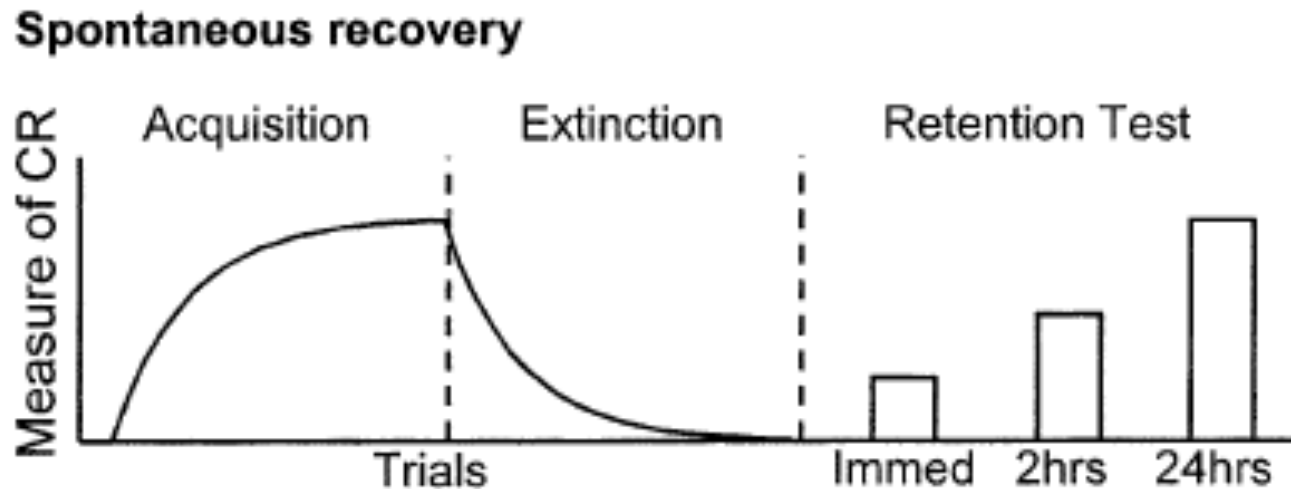


Hippocampal lesions handicap the model's ability to infer new clusters

Why are memories hard to modify?

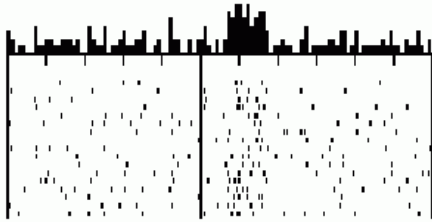


Relapse in classical conditioning

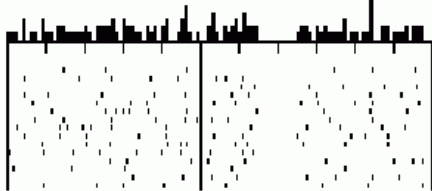


Prediction errors and learning

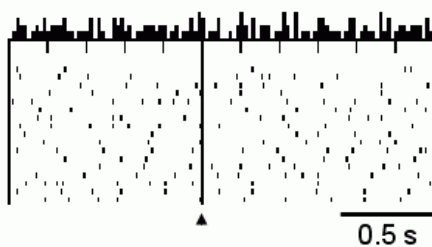
Conditioned stimulus
predicting reward



Conditioned stimulus
predicting absence of reward

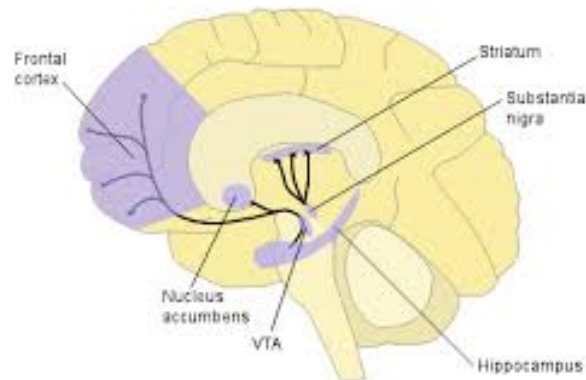


Known neutral stimulus



Rescorla-Wagner model

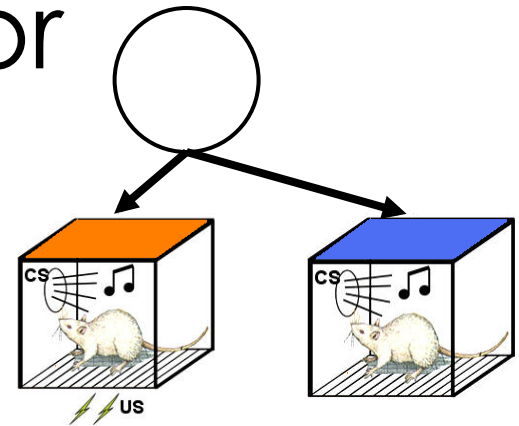
$$V \leftarrow V + \eta \underbrace{[\text{outcome} - V]}_{\text{prediction error}}$$



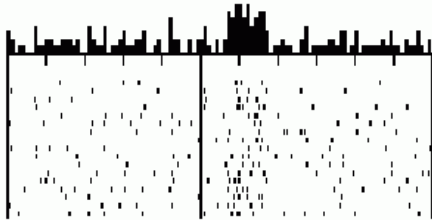
Schultz (1998)

An alternative view: two roles for prediction error

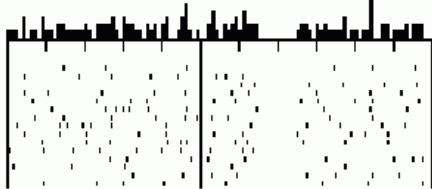
When errors are small:
memory modification



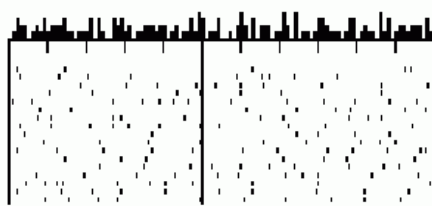
Conditioned stimulus
predicting reward



Conditioned stimulus
predicting absence of reward

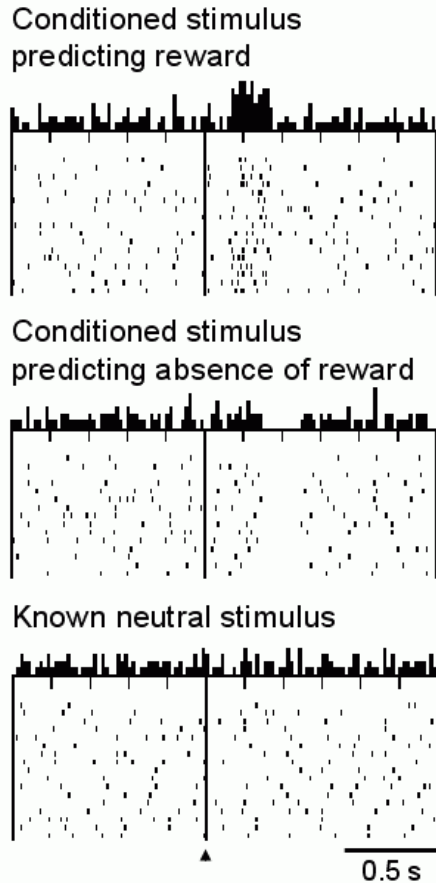


Known neutral stimulus



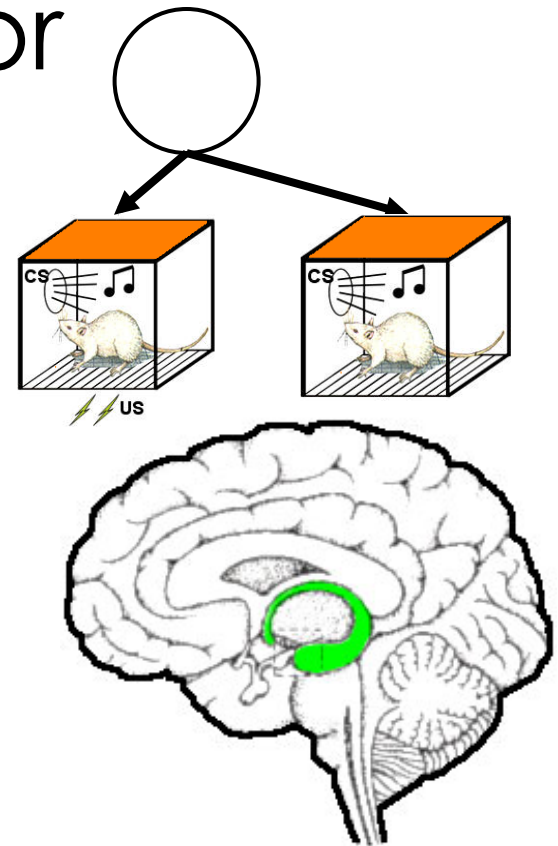
0.5 s

An alternative view: two roles for prediction error

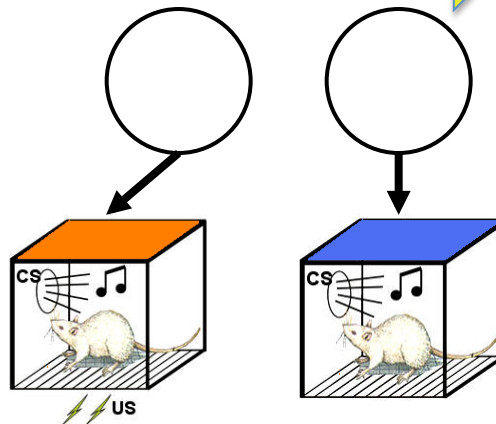


Schultz (1998)

When errors are small:
memory modification



When errors are large:
memory formation



How to erase a fear memory

- Prediction errors should be large enough to drive learning, but not so large that a new latent cause is inferred.

How to erase a fear memory

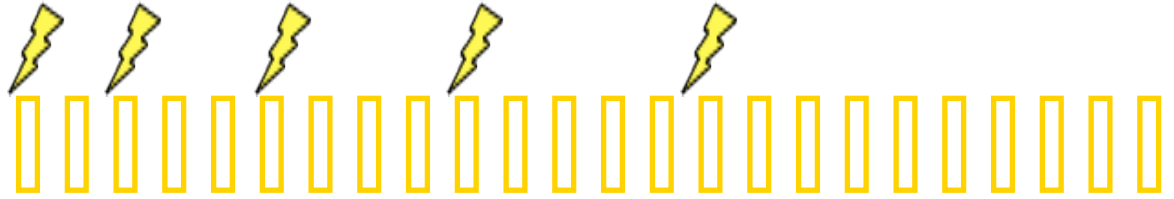
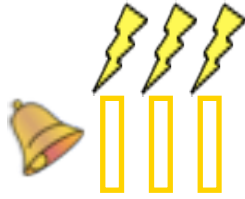
- Prediction errors should be large enough to drive learning, but not so large that a new latent cause is inferred.
- Titrate prediction errors by **extinguishing gradually**.

Testing the model: gradual extinction

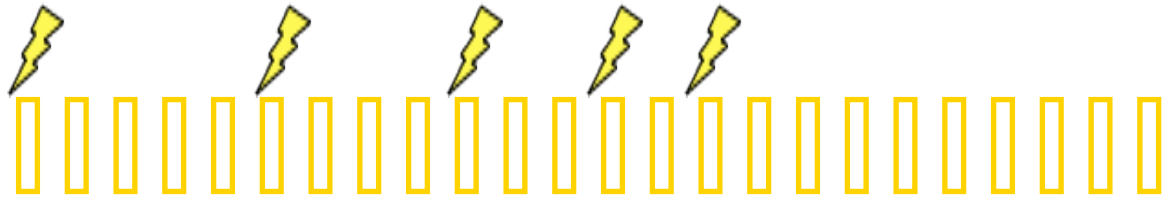
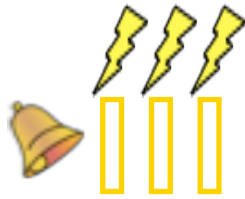
acquisition

extinction

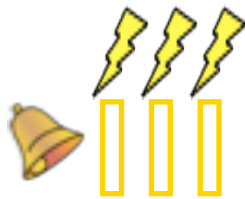
gradual extinction

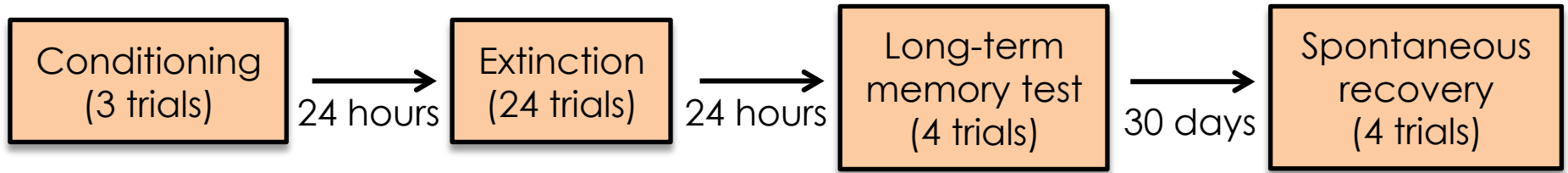


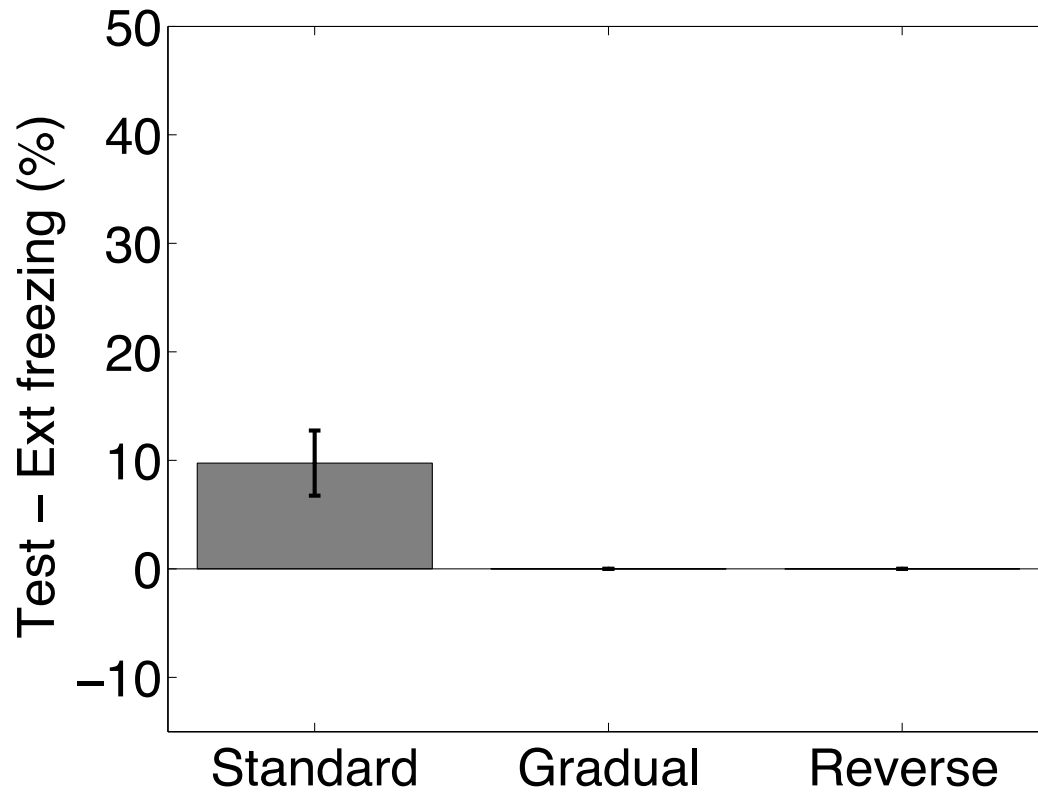
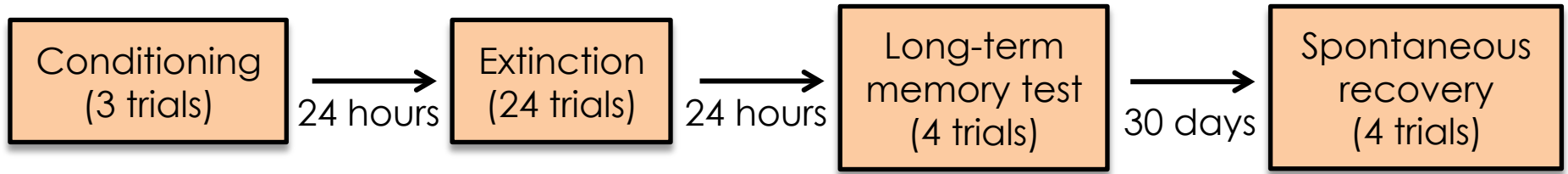
gradual reverse

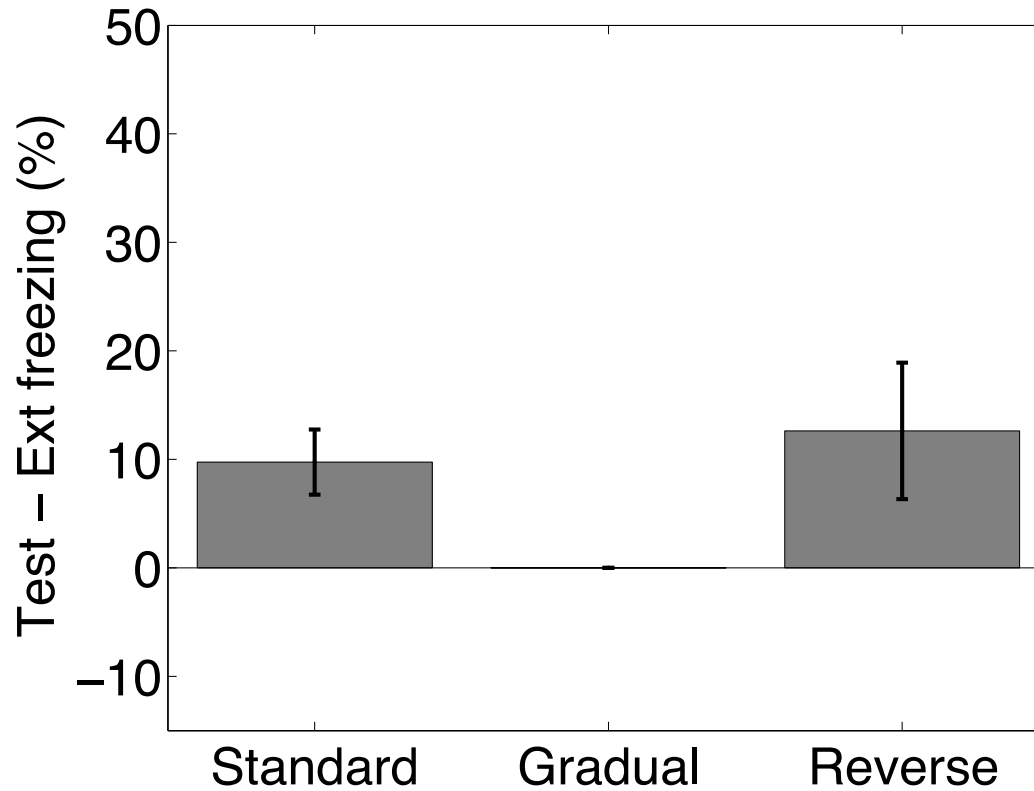
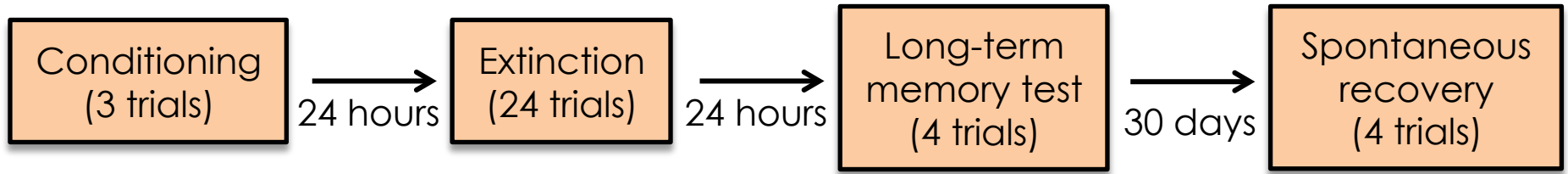


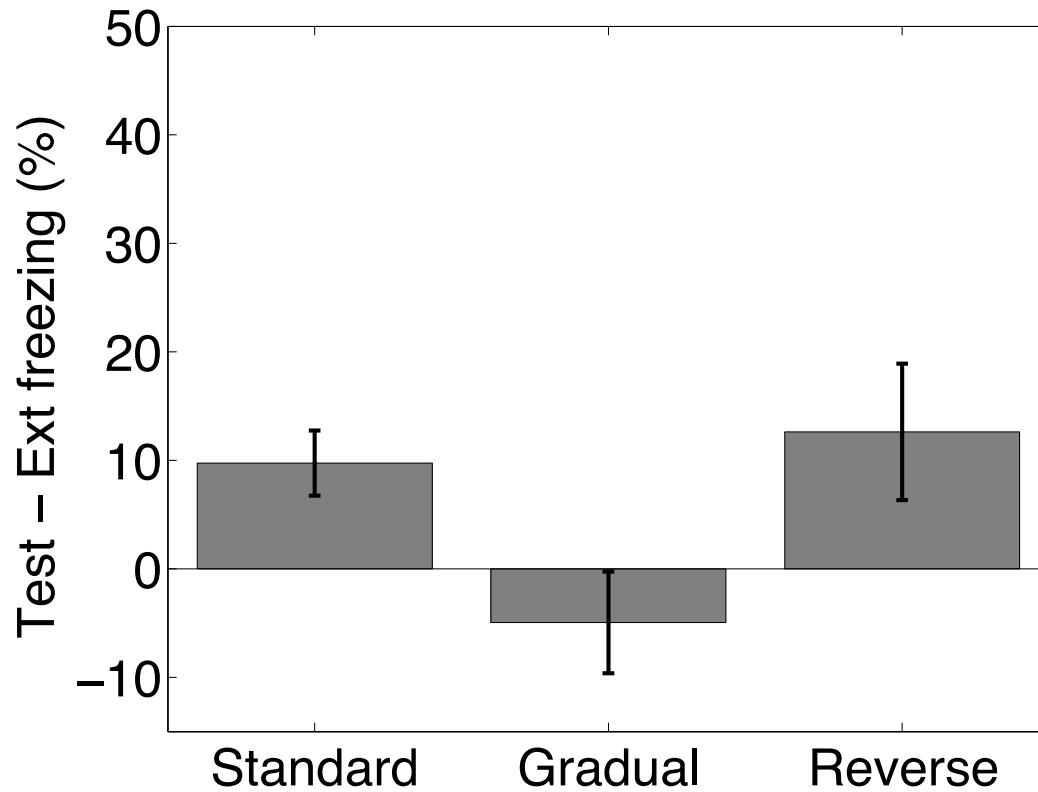
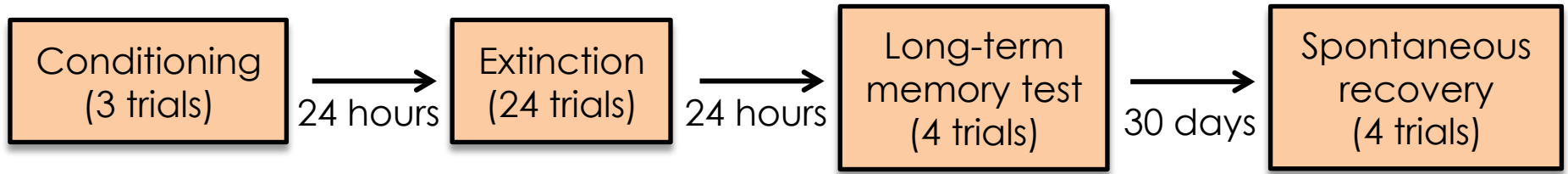
regular extinction



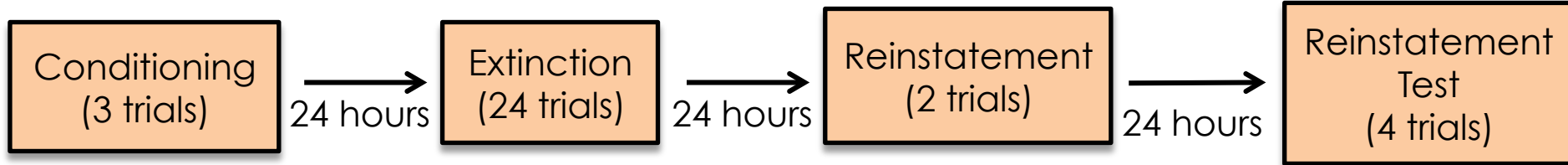




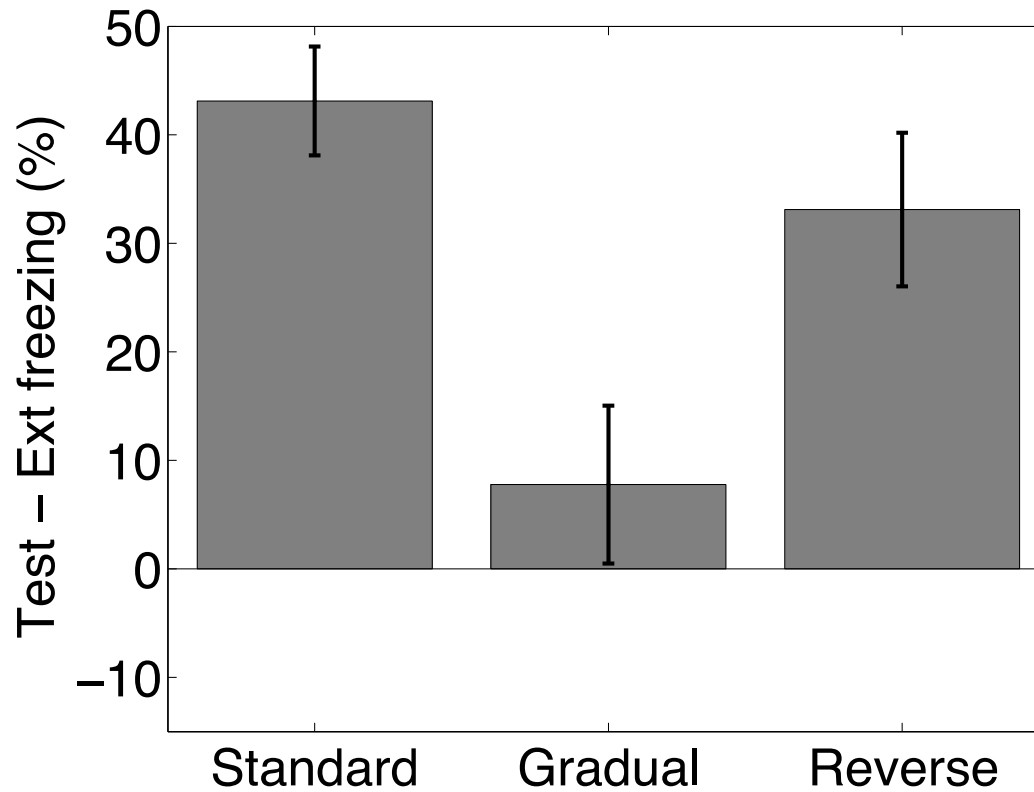
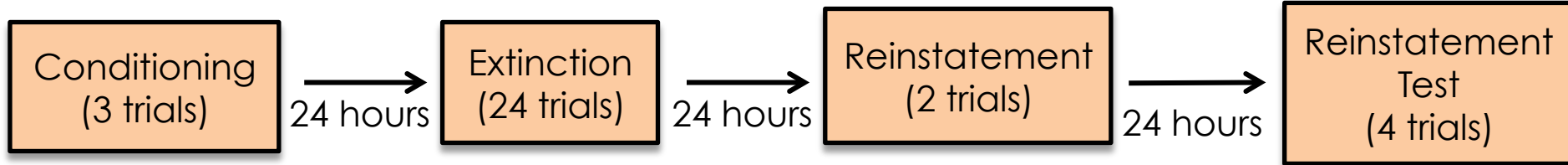




Reinstatement design

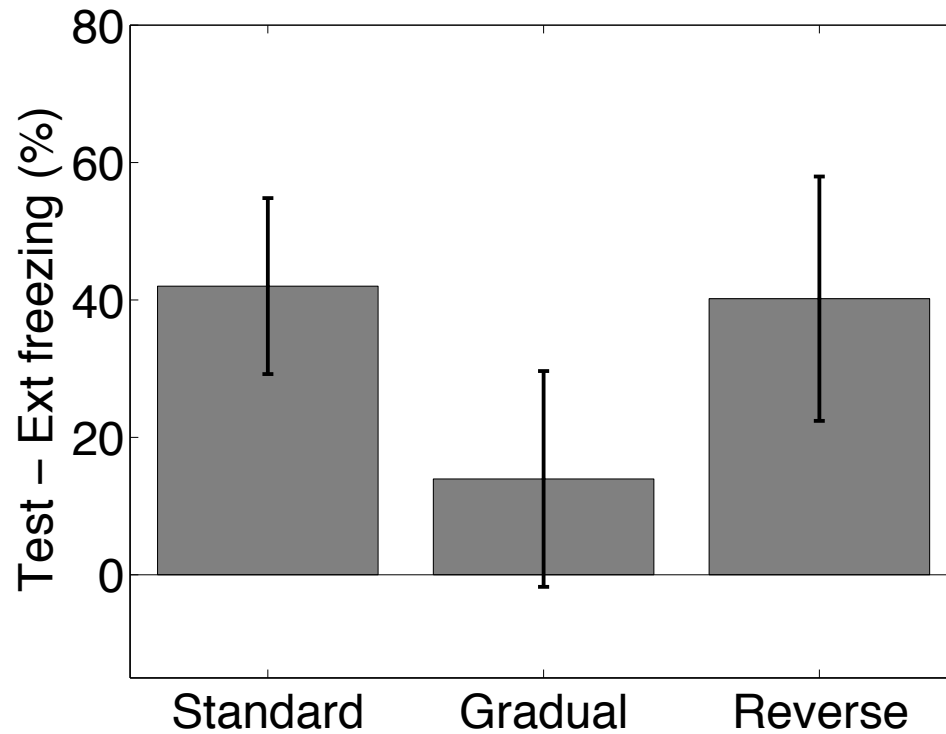


Reinstatement design



Gradual extinction in humans

Gradual extinction in humans



Predicting spontaneous recovery in humans



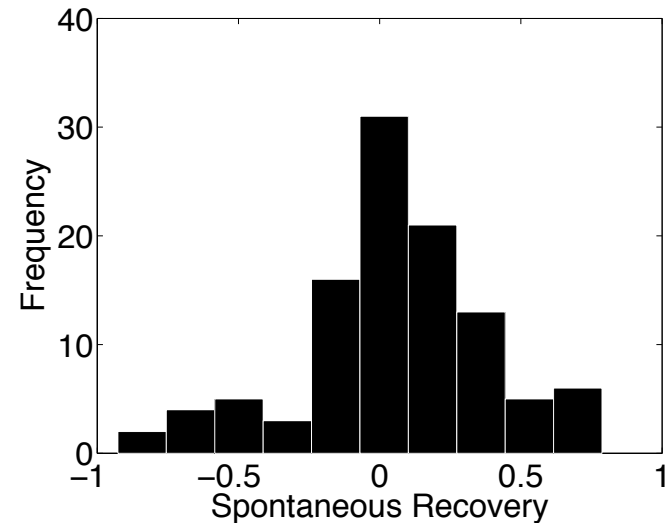
(partial reinforcement)

Predicting spontaneous recovery in humans



(partial reinforcement)

Why do some people show a return of fear, and some don't?

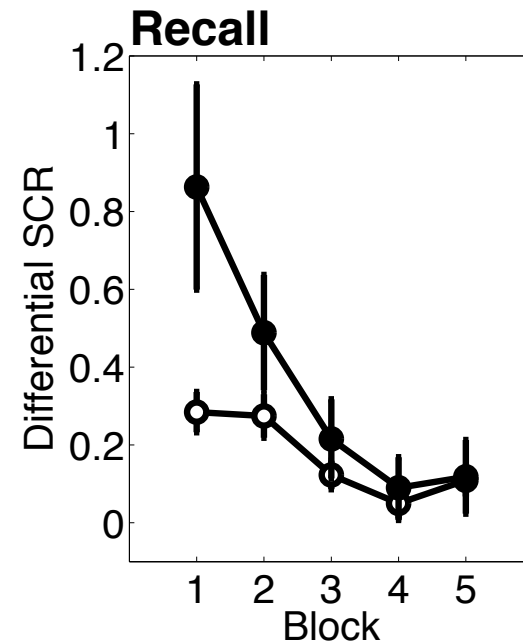
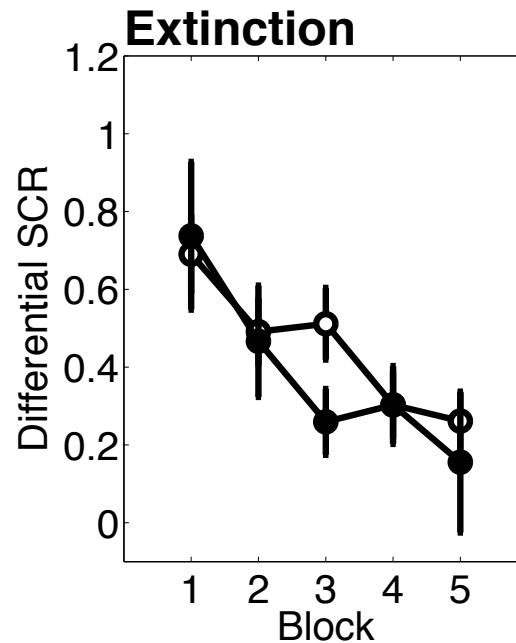
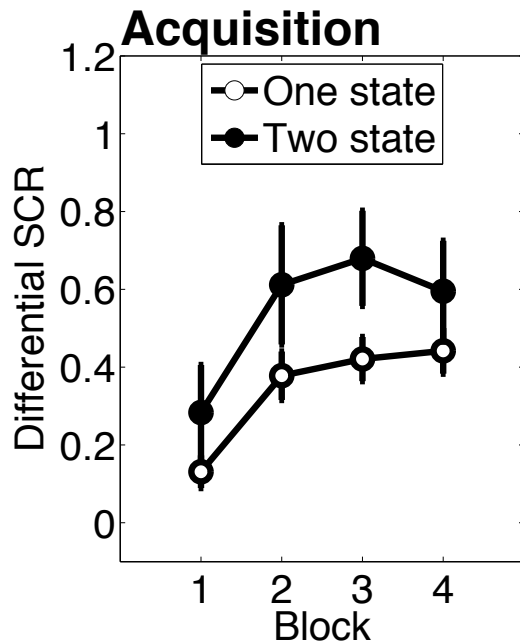


Gershman & Hartley
(submitted)

Predicting spontaneous recovery in humans



(partial reinforcement)

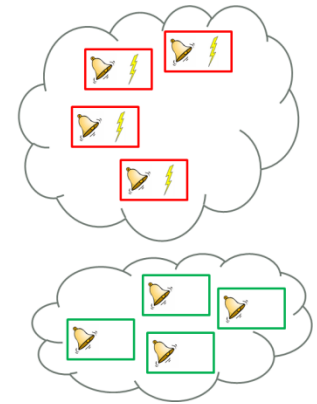
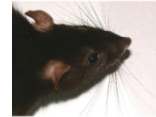


Model, fit only to conditioning & extinction data, divides subjects into two groups

Gershman & Hartley (submitted)

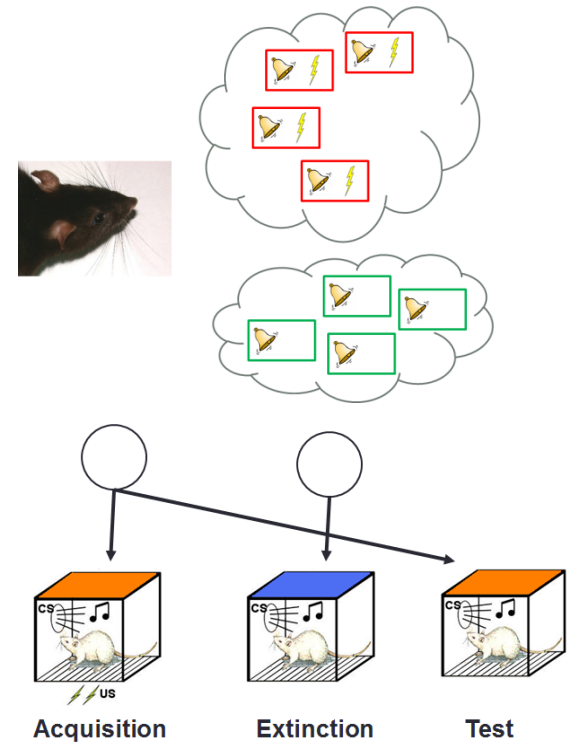
Summary

- Conditioning as clustering



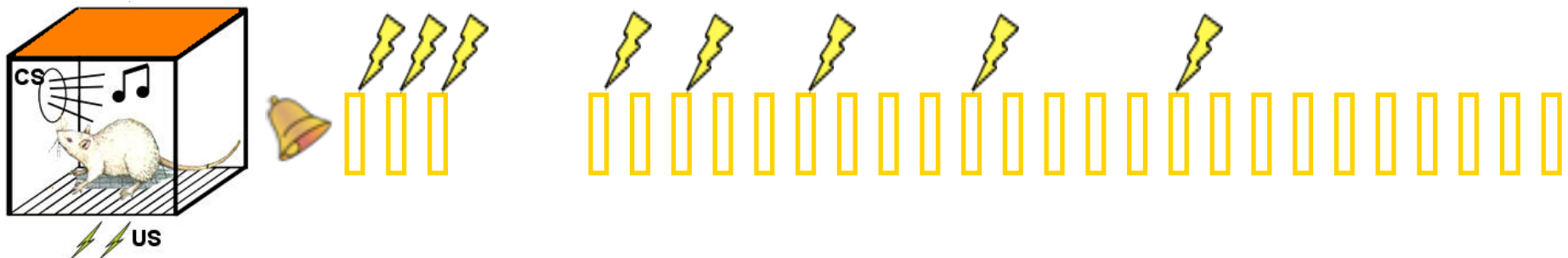
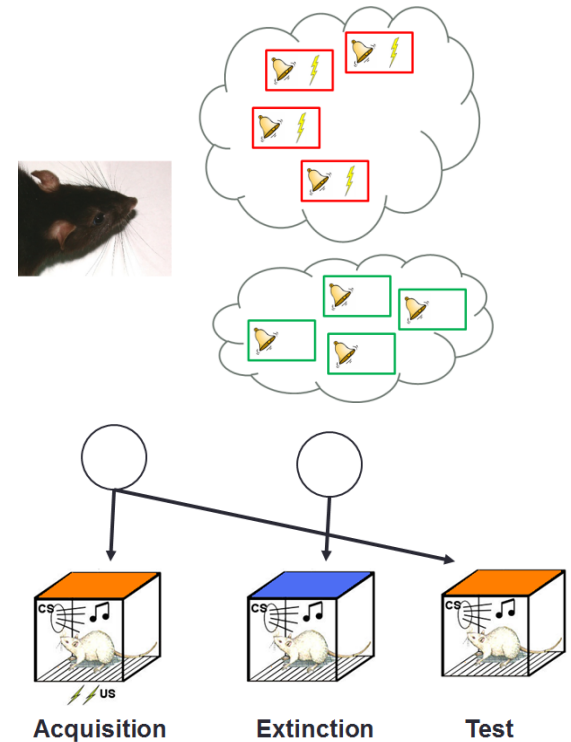
Summary

- Conditioning as clustering
- Memories reflect inferences about latent causes



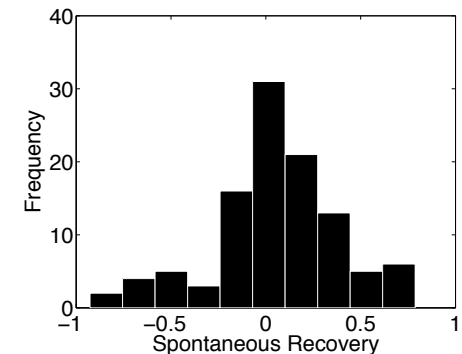
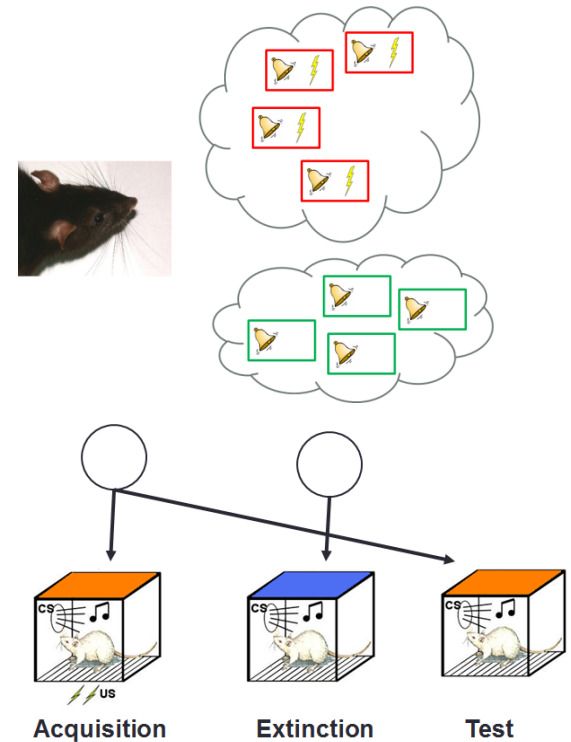
Summary

- Conditioning as clustering
- Memories reflect inferences about latent causes
- Gradual extinction prevents the return of fear



Summary

- Conditioning as clustering
- Memories reflect inferences about latent causes
- Gradual extinction prevents the return of fear
- Explaining individual differences in the return of fear



Big picture

- When do we modify old memories, and when do we create new ones?
- This question can be answered within a probabilistic computational framework:
we create new memories when we infer new latent causes in our environment
- This principle has deep explanatory power across multiple domains

Acknowledgments



Ken Norman
(Princeton)



Yael Niv
(Princeton)



David Blei
(Princeton)



Marie Monfils
(UT Austin)

Carolyn Jones (UT Austin)
Cate Hartley (Sackler Institute)
Liz Phelps (NYU)