



Princeton Computational Memory Lab

Testing Psychological Theories with Multivariate Pattern Analysis

Ken Norman

Department of Psychology and Neuroscience Institute

Princeton University

July 25, 2008

Uses of Multivariate Pattern Analysis (MVPA)

- Determine **where** and **how** information is represented in the brain
- Track time-varying cognitive states
 - Multivariate methods improve sensitivity
 - One benefit of this extra sensitivity is that we can generate a useful estimate of the subject's cognitive state at a **particular point in time**
 - We can leverage this temporal sensitivity to test psychological theories

Theory testing

- Psychological theories can be viewed as if-then statements
- If [COGNITIVE STATE] then [OUTCOME]
- Standard, behavioral approach to testing theories:
 - Set up experimental conditions that you hope will bring about the cognitive state of interest
 - Look for the predicted outcome
- Problem: Our ability (as experimenters) to control the subject's cognitive state is limited
- If you don't get the effect that you want, it may be because the theory is wrong, or it may be that you weren't successful in eliciting the cognitive state of interest

Theory testing

- Our approach: Use pattern classification algorithms, applied to brain imaging data, to isolate **distributed patterns of neural activity** corresponding to cognitive states of interest
- Once we have trained a pattern classifier to detect the neural correlate of a particular cognitive state, we can use the classifier to track fluctuations in that cognitive state over time
- We can use this time-varying readout of the subject's cognitive state to test hypotheses about how that cognitive state relates to behavior
- Another benefit: We can use more open-ended, naturalistic experimental designs

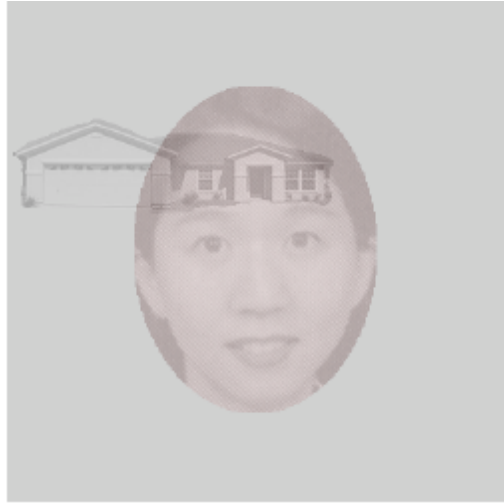
Outline

- Case studies
 - Using classifiers (applied to EEG) to test a theory of how brain activity drives learning
 - Using classifiers (applied to fMRI) to track cognitive processes during memory search
- Limitations of the classifier approach

General Design

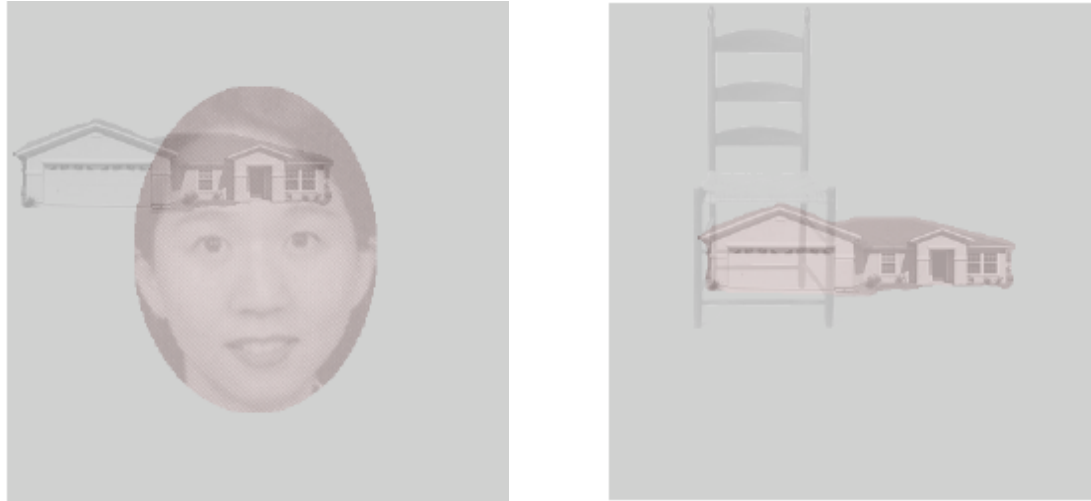
- All of the experiments that I will present have the same 2-part design
- Part 1: Collect data for classifier training
- Strongly & unambiguously elicit the cognitive states of interest. Use these data to train the classifier
- Part 2: Generalization
- Apply the trained classifier to situations where the subject's cognitive state is more variable
- Use the classifier's readout of the subject's cognitive state to predict behavior

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)



- Key finding: Subjects are faster to respond to stimuli that were previously **attended** and slower to respond to stimuli that were previously **ignored** (e.g., Tipper, 1985)

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)



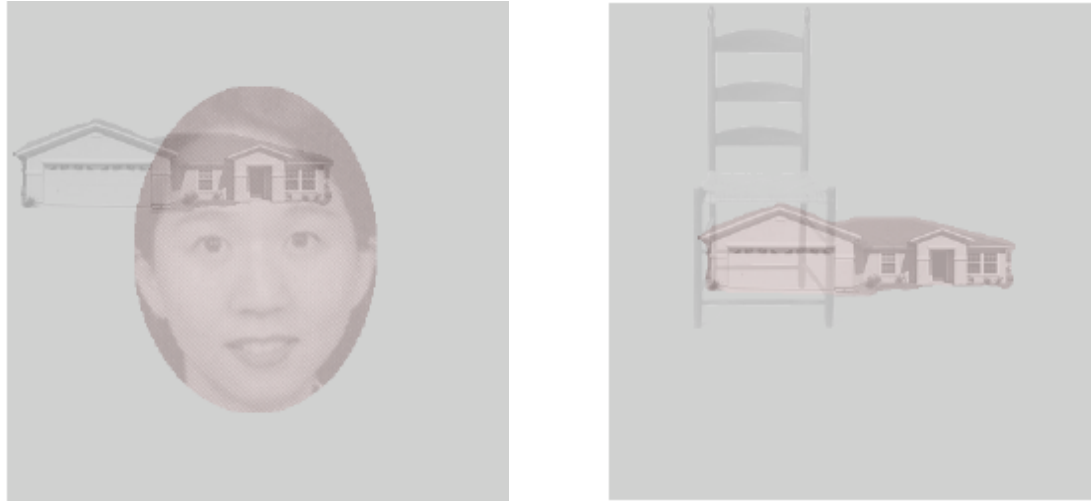
- Theory: competition drives learning (Norman et al., 2007).
When two representations compete...

Winning representation ➡ strengthened

Representations that **activate but lose** ➡ weakened

Representations that **do not activate** ➡ no change

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)



- Theory: competition drives learning (Norman et al., 2007).
When two representations compete...

targets (red)	→ win	→ strengthened
distractors (gray)	→ activate but lose	→ weakened
others	→ do not activate	→ no change

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)

targets (red)	➡ win	➡ strengthened
distractors (gray)	➡ activate but lose	➡ weakened
others	➡ do not activate	➡ no change

- Problem: Negative priming effects are not always found
- Explanation 1: Theory is wrong => representations that lose the competition are not weakened

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)

targets (red)	→ win	→ strengthened
distractors (gray)	→ activate but lose	→ weakened
others	→ do not activate	→ no change

- Problem: Negative priming effects are not always found
- Explanation 1: Theory is wrong => representations that lose the competition are not weakened

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)

targets (red)	→ win	→ strengthened
distractors (gray)	→ activate but lose	→ weakened
others	→ do not activate	→ no change

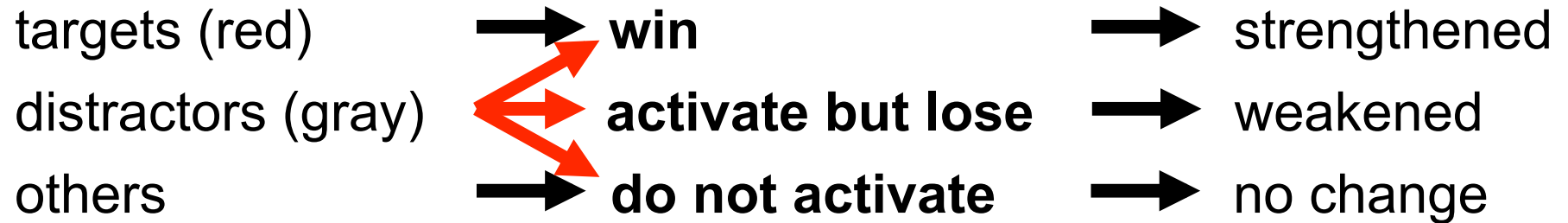
- Problem: Negative priming effects are not always found
- Explanation 1: Theory is wrong => representations that lose the competition are not weakened
- Explanation 2: Mapping between experimental conditions and activation dynamics is noisy

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)

targets (red)	→ win	→ strengthened
distractors (gray)	→ activate but lose	→ weakened
others	→ do not activate	→ no change

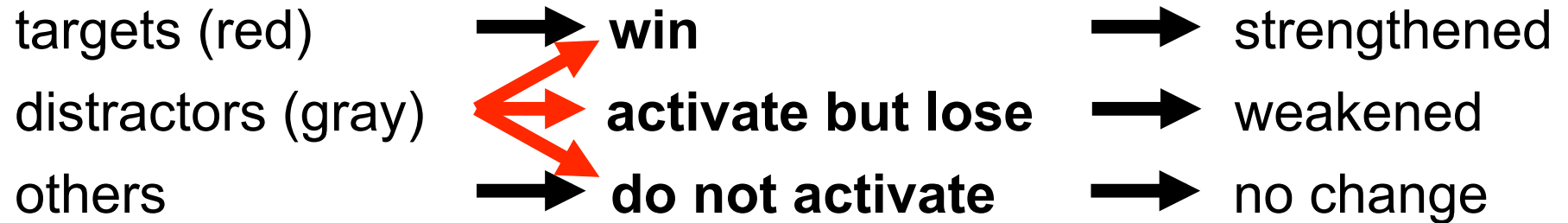
- Problem: Negative priming effects are not always found
- Explanation 1: Theory is wrong => representations that lose the competition are not weakened
- Explanation 2: Mapping between experimental conditions and activation dynamics is noisy

Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)



- Problem: Negative priming effects are not always found
- Explanation 1: Theory is wrong => representations that lose the competition are not weakened
- Explanation 2: Mapping between experimental conditions and activation dynamics is noisy

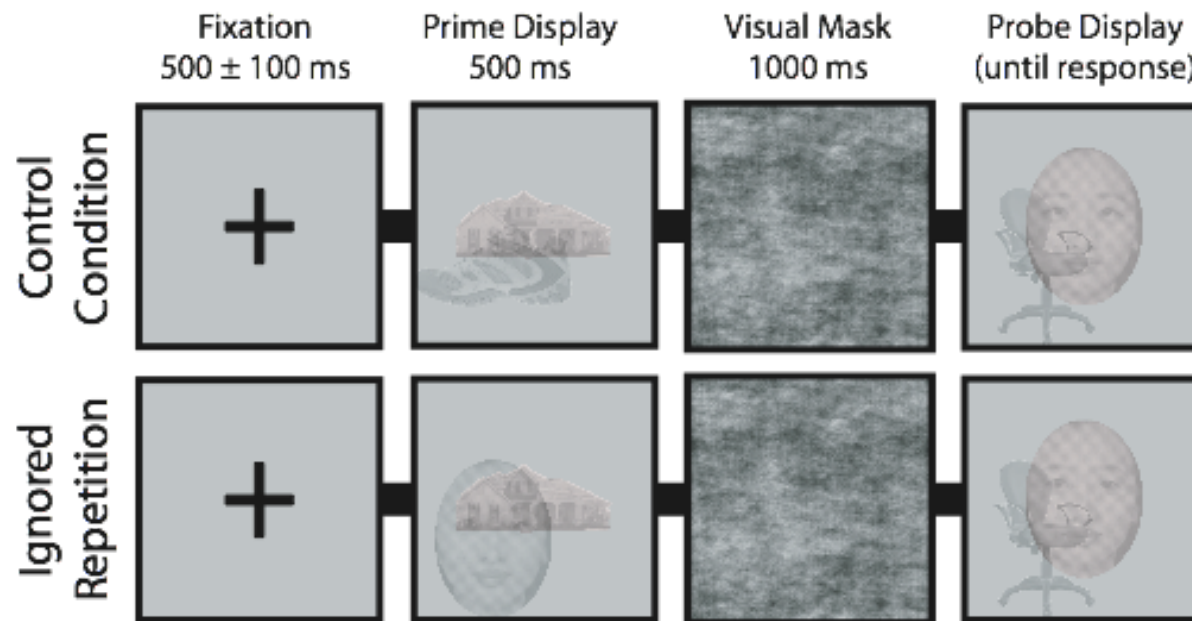
Case Study: Negative Priming and Competitor Weakening (Newman & Norman, in prep)



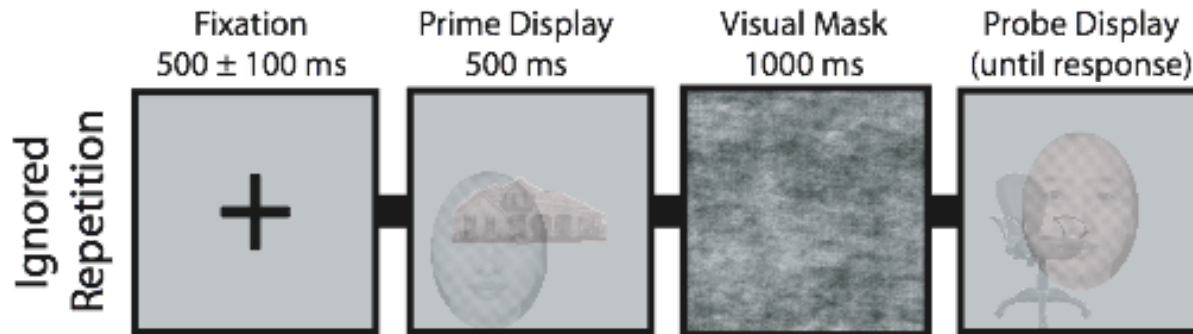
- Our approach:
- Use pattern classifiers to **directly measure** distractor activation
- Use this measure of distractor activation to predict whether subjects will show **positive priming**, **negative priming**, or **no priming** for that item

Delayed Match to Sample Paradigm

- stimuli were composed of a red-tinted **target** stimulus on top of a black & white **distractor**
- stimuli were either faces, houses, shoes, or chairs
- targets and distractors were always from different categories

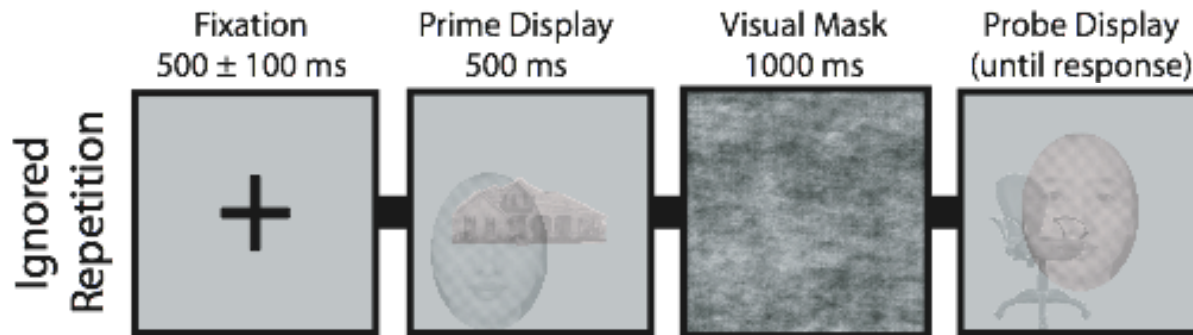


Analysis Strategy



- Part 1: Apply classifiers to patterns of EEG data collected during the prime; train classifiers to read out the category of the **prime target** stimulus
- Part 2: Use these trained classifiers to read out how much subjects were processing the **prime distractor** stimulus
- Use readout of distractor activity to predict RT to the probe
- Training & testing were always done on different parts of the data set

Classification Details



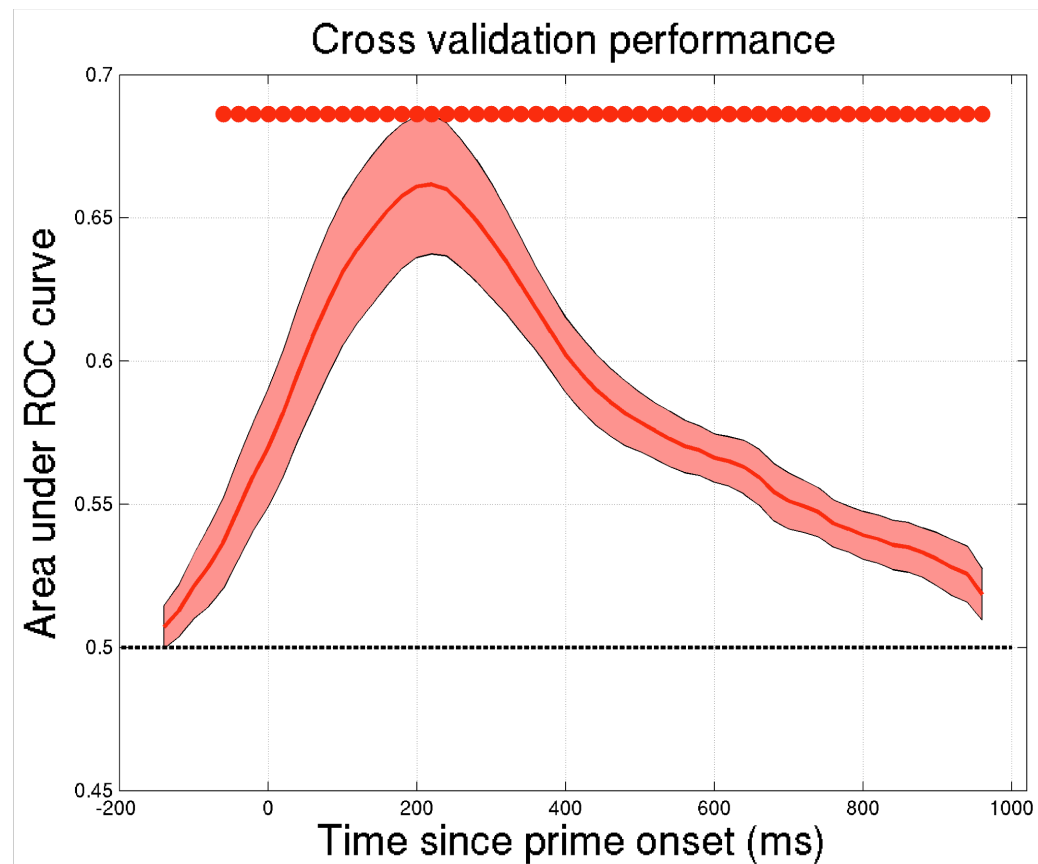
- We applied classifiers to EEG data from the 1000 ms window starting with prime onset
- Separate classifiers were used for each category
 - Face-on-screen-as-target vs. face-absent
 - Shoe-on-screen-as-target vs. shoe-absent
 - Chair-on-screen-as-target vs. chair-absent
 - House-on-screen-as-target vs. house-absent

Classification Details

- Record EEG from 77 electrodes
- EEG time series were spectrally decomposed into wavelet power coefficients at 49 frequencies (ranging from 2 to 128Hz)
- Wavelet time-series were down-sampled into 20ms time bins (50 time bins per 1000ms trial)
- For classification purposes, our features were defined by the crossing of [77 electrodes] X [49 frequencies] X [50 time bins]
- We discarded features that did not (individually) discriminate between the conditions of interest in the training set
- Ridge regression classifier
- We trained a separate classifier for each time bin

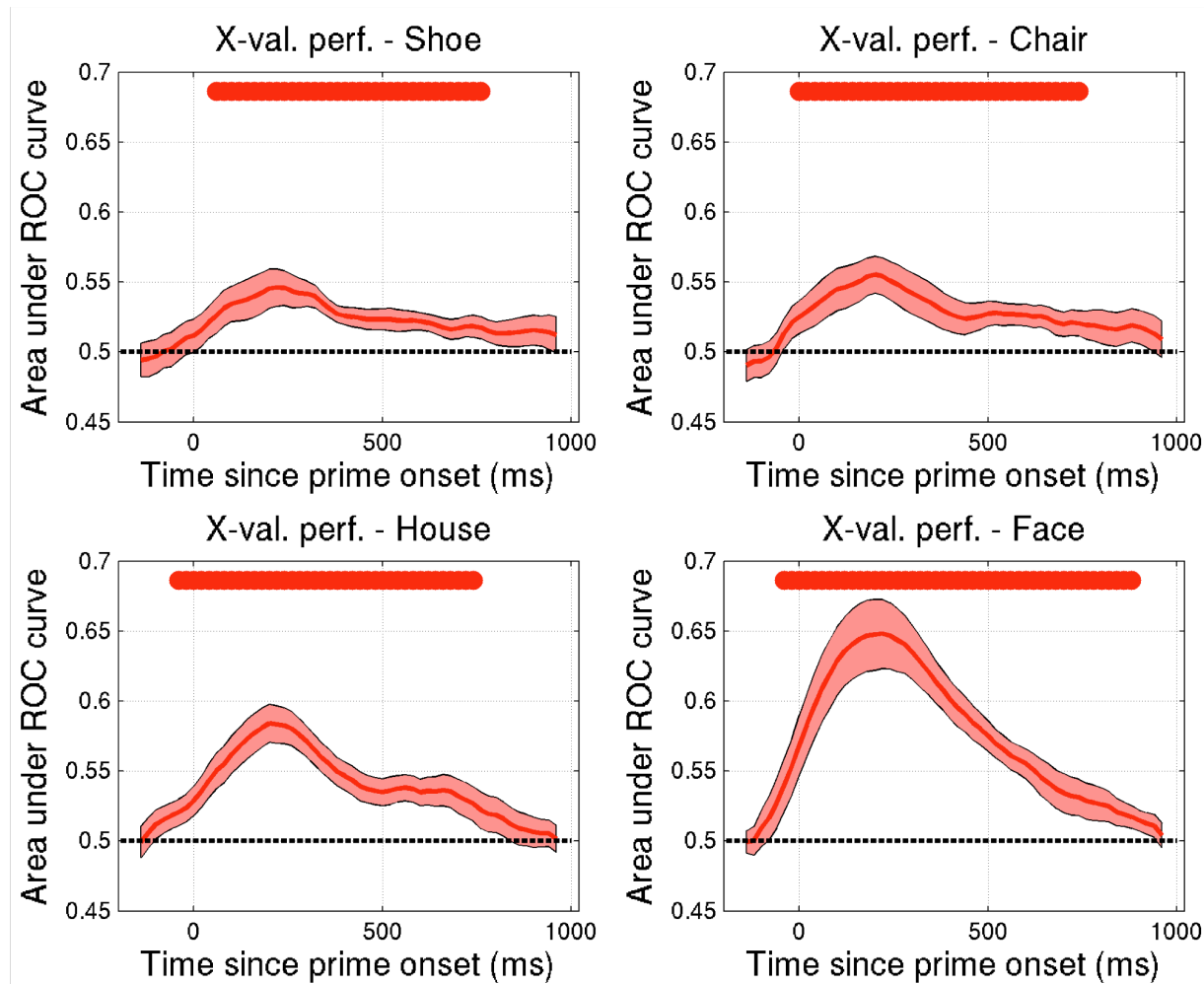
Part 1 Results: Target Classification

- average of all 4 categories



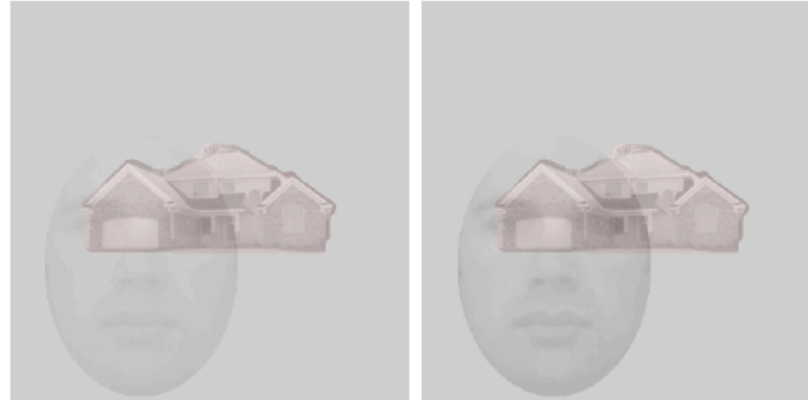
Part 1 Results: Target Classification

- target classification was above chance for all 4 categories



Part 2: Classifying Distractors

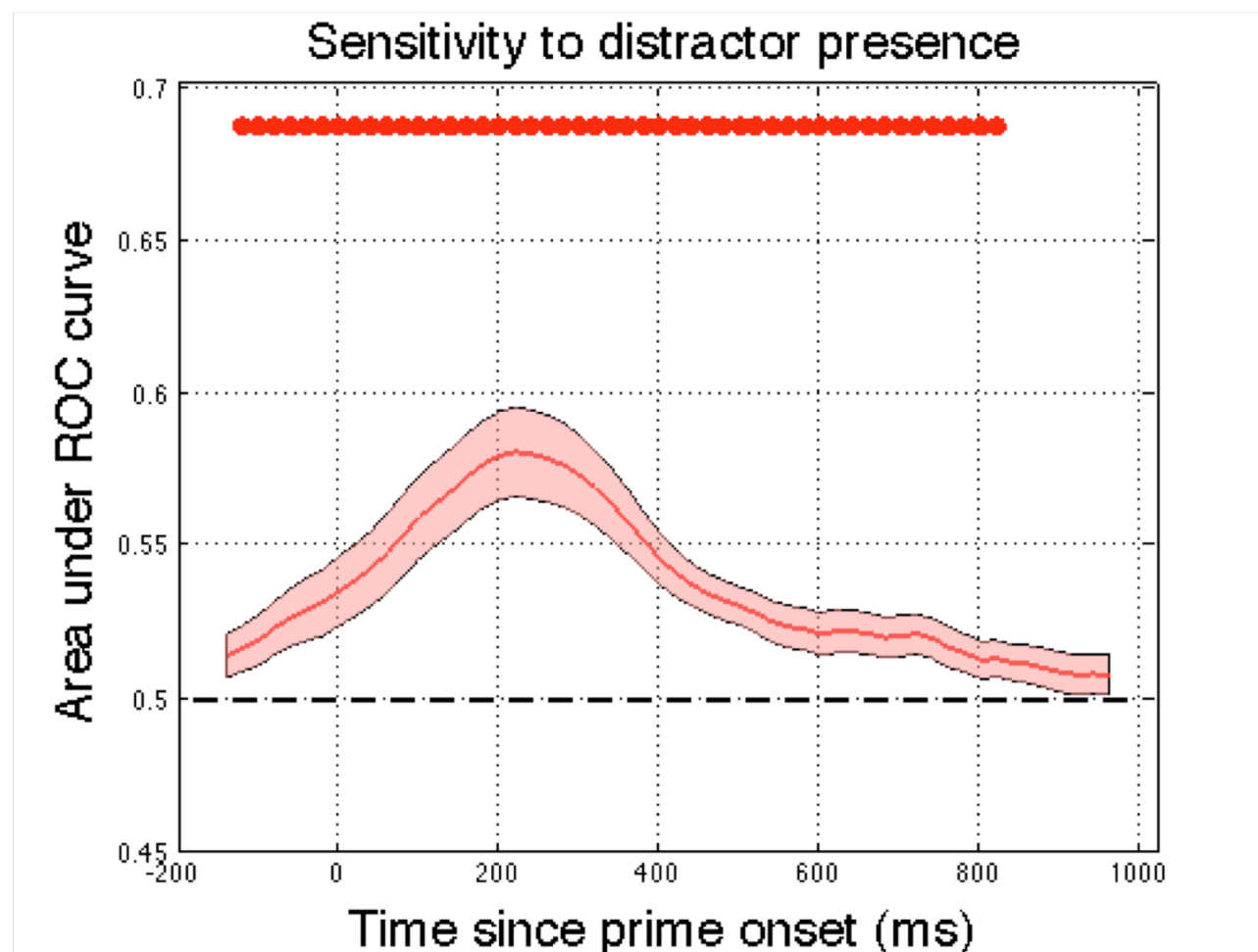
- Key question: Can we use the classifiers that were trained to detect the **target** category to also detect the **distractor** category?
- Compute average classifier activity when a category is the **distractor** vs. **not present**
- Vary distractor strength

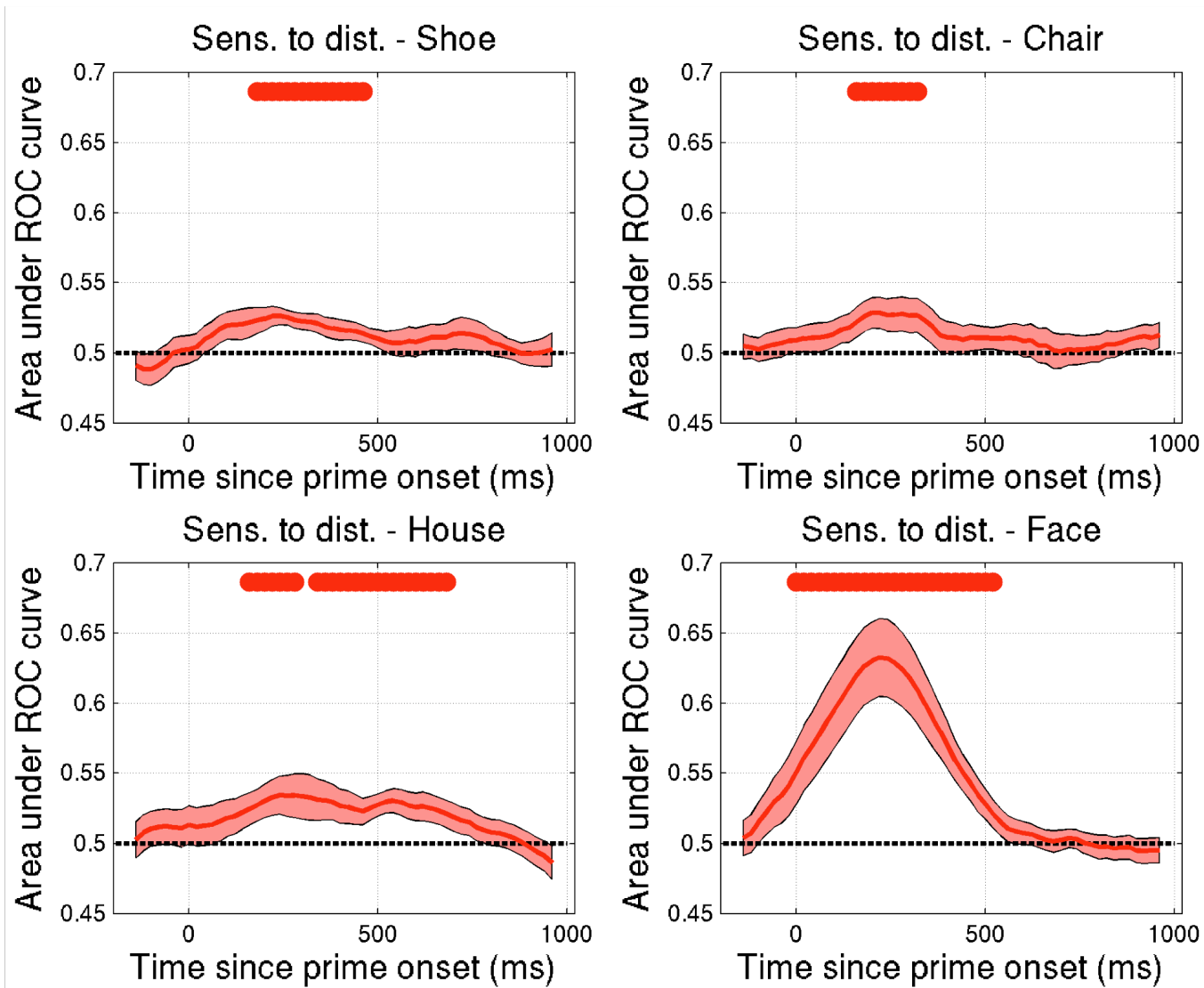


Weak
distractor
cue

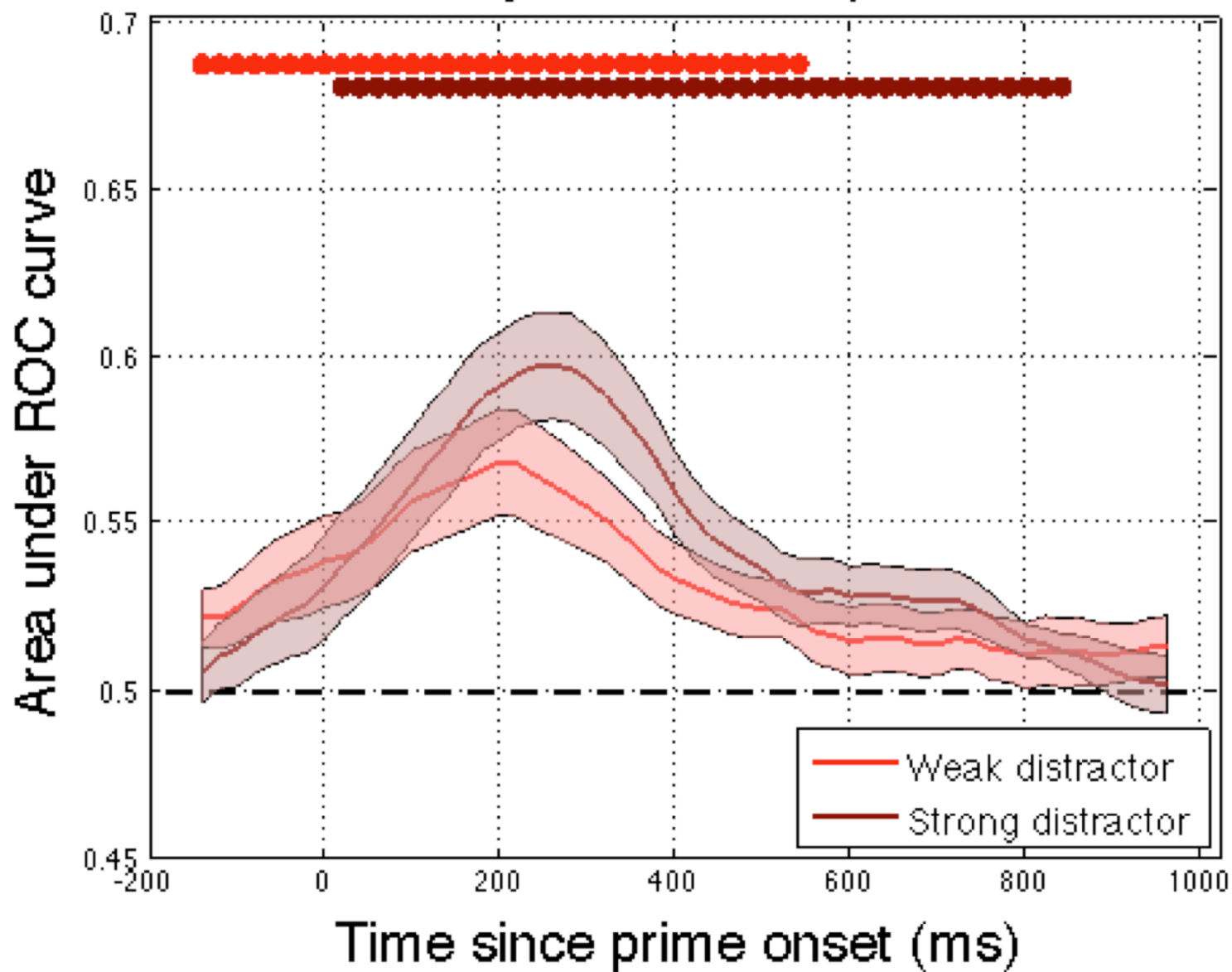
Strong
distractor
cue

- The classifier's readout of distractor activity should be higher in the **strong distractor** condition than the **weak distractor** condition

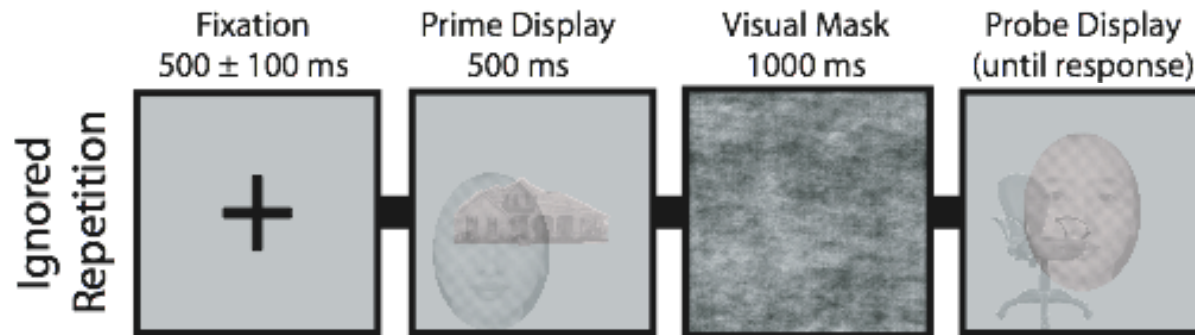




Sensitivity to distractor presence



Part 2: Relating Distractor Activity to RT



- Predictions of the competition-dependent learning theory:

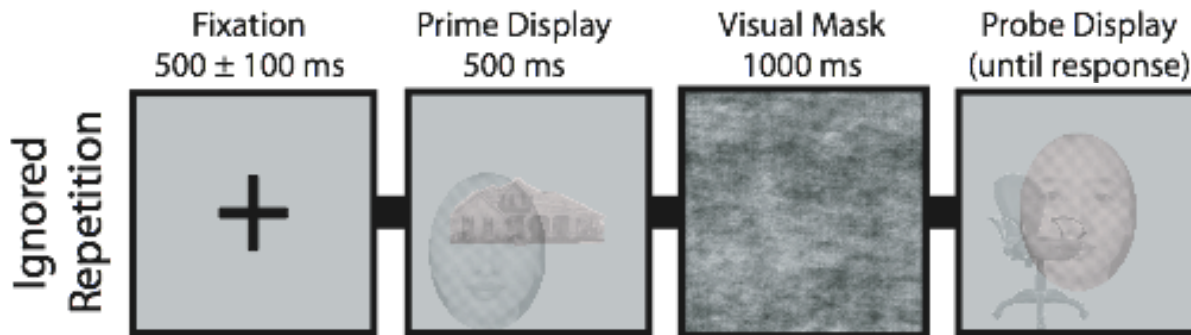
- If the distractor:

does not activate → no change

activates but loses → weakened (neg. priming)

wins → strengthened (pos. priming)

Part 2: Relating Distractor Activity to RT



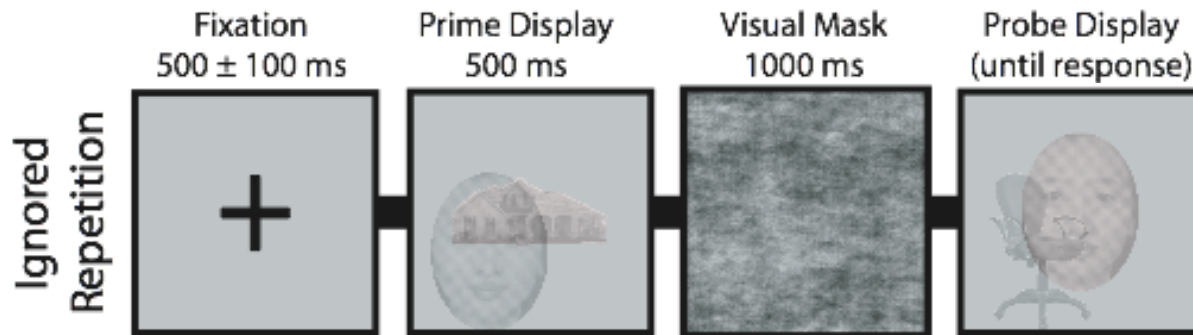
- Predictions of the competition-dependent learning theory:

- If the distractor:

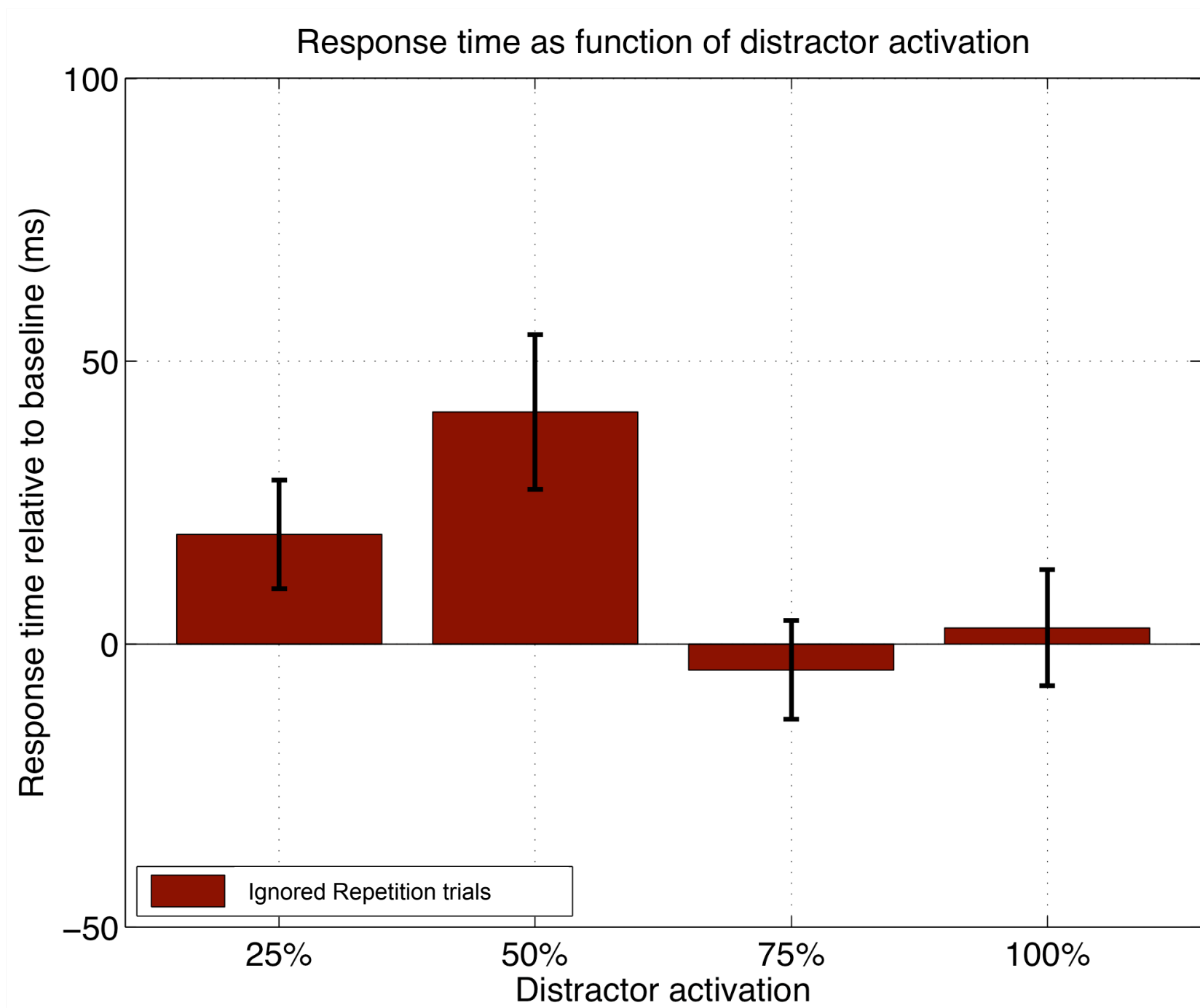
LOW	does not activate	→	no change
MED	activates but loses	→	weakened (neg. priming)
HIGH	wins	→	strengthened (pos. priming)

- Key prediction: negative priming effect should be largest for **moderate** levels of distractor activity

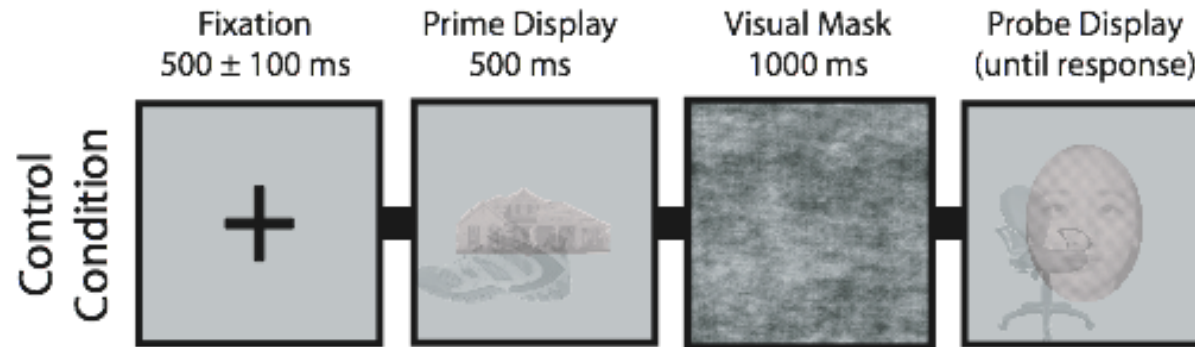
Part 2: Relating Distractor Activity to RT



- Key prediction: negative priming effect should be largest for **moderate** levels of distractor activity
- To test this, we split trials into quartiles based on distractor activity (averaged across time bins) and computed the priming effect as a function of distractor activity

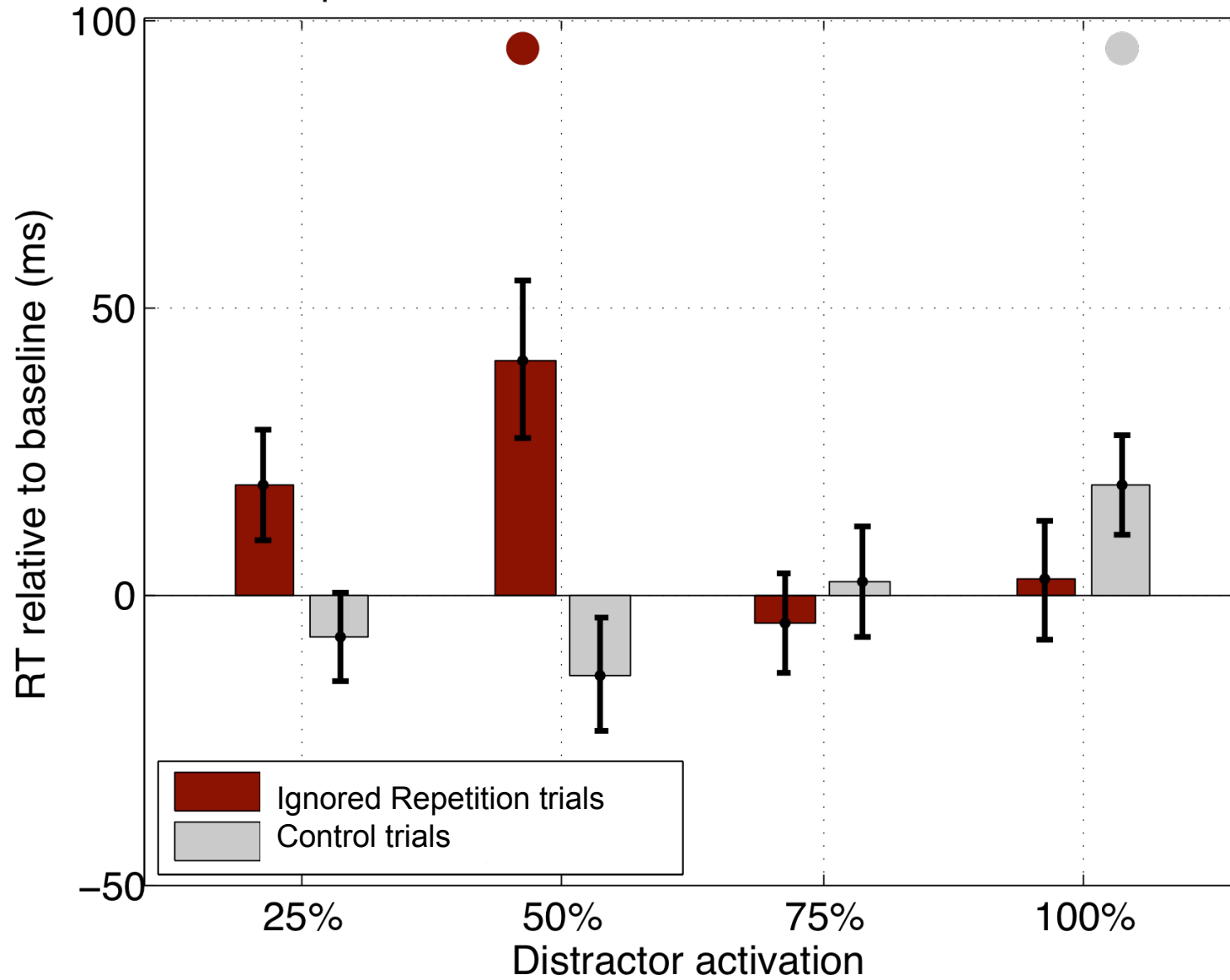


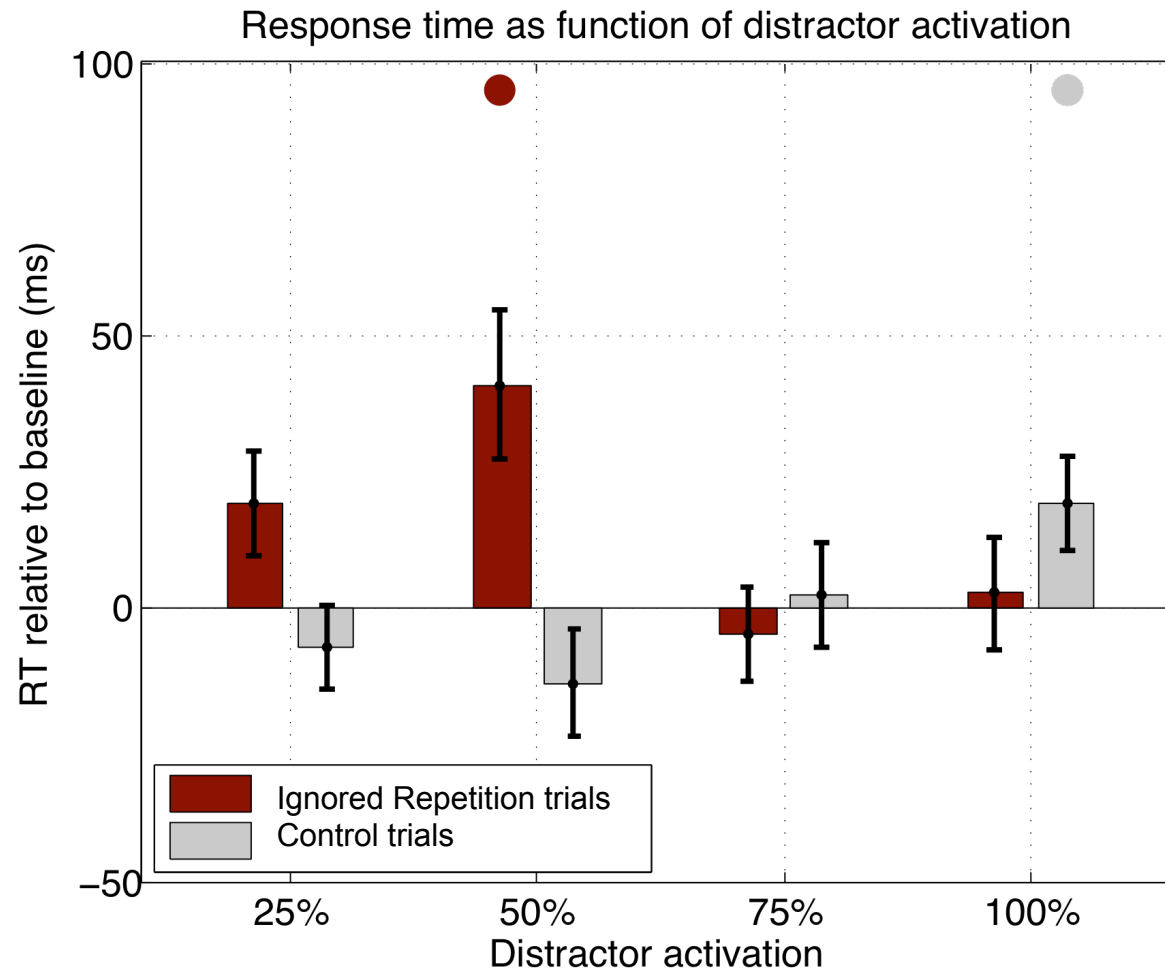
Control Condition Results



- To assess whether these effects were specifically due to priming, we also ran this analysis on control trials
- In control trials, the prime and probe use completely different categories

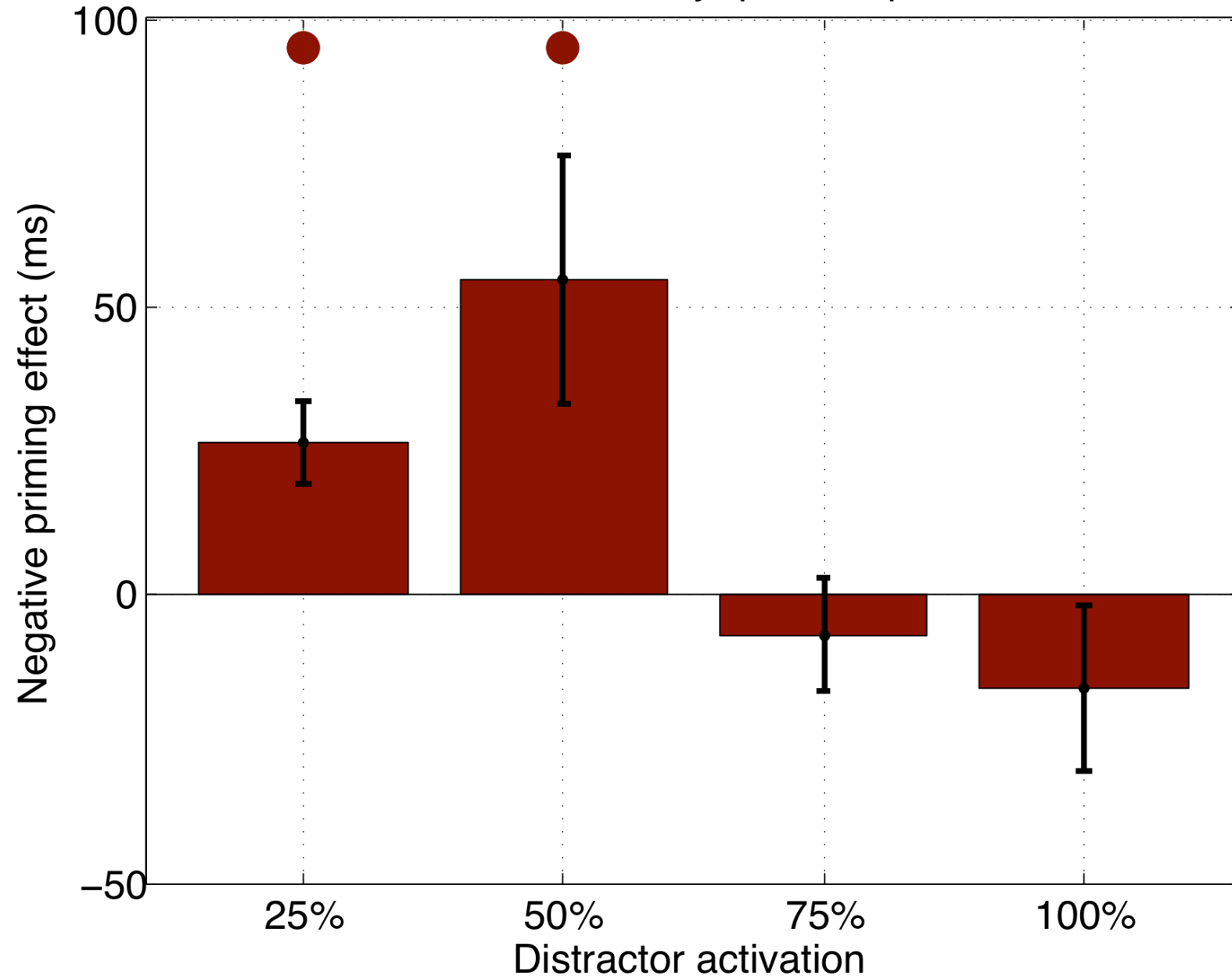
Response time as function of distractor activation





- General attentional effect:
- IF subjects are not focusing on the target, then we should see:
 - High distractor activation during the prime
 - Slow responding to the target during the probe

Matched NP by quartile split



Negative Priming: Summary

- The results from this study support our hypothesis that **moderate** activation of a neural representation leads to weakening of that representation
- Sorting trials by distractor activation allowed us to isolate a large, robust negative priming effect
 - NP effect across all trials (no sorting) = 14 ms
 - NP effect given moderate activation = 51 ms
- These findings fit well with results from studies of LTD
 - At the synaptic level, moderate depolarization of the postsynaptic neuron leads to LTD (e.g., Artola et al., 1990)
 - Our study demonstrates this dynamic at the level of human behavior

Negative Priming: Summary

- We were able to leverage highly-classifiable cognitive states as a “contrast dye” to improve temporal resolution
 - My lab has no intrinsic interest in faces, houses, shoes, & chairs
 - We used the categorical stimuli in the priming study because they are highly classifiable, and this allowed us to derive a useful trial-by-trial measure of distractor processing

Case Study 2: Tracking Memory Search

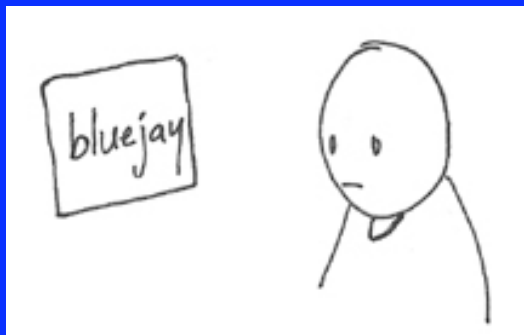
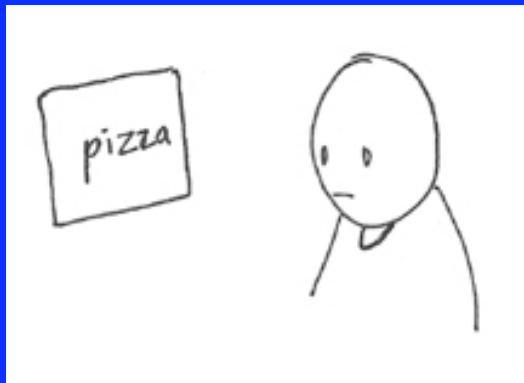
- How do we search memory for a particular event?
- Web search analogy:
 - To find the web page you're looking for, you need to use the right search terms
 - The same thing applies to memories
 - A huge portion of the variance in what you retrieve depends on how you cue memory
- The goal of the work I am going to describe is to **read out the information that subjects are using to search memory**, and to relate it to their behavior

Case Study 2: Tracking Memory Search

- Quick overview of the task that we use to study memory search, and theories of memory search
- Two fMRI studies of free recall

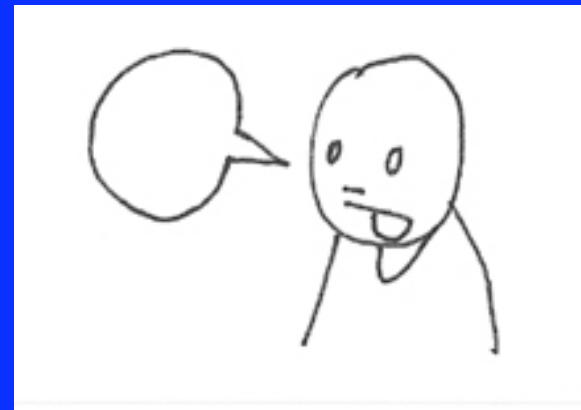
Memory search in the lab

- Free recall

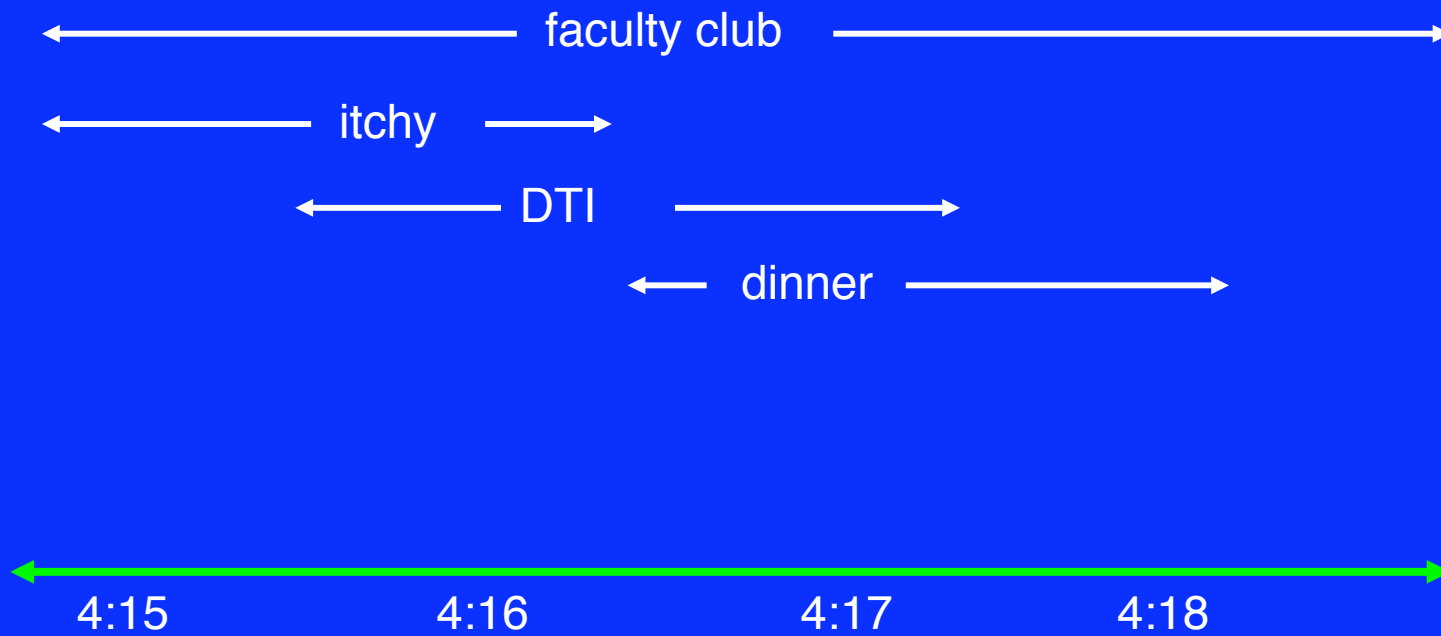


etc...

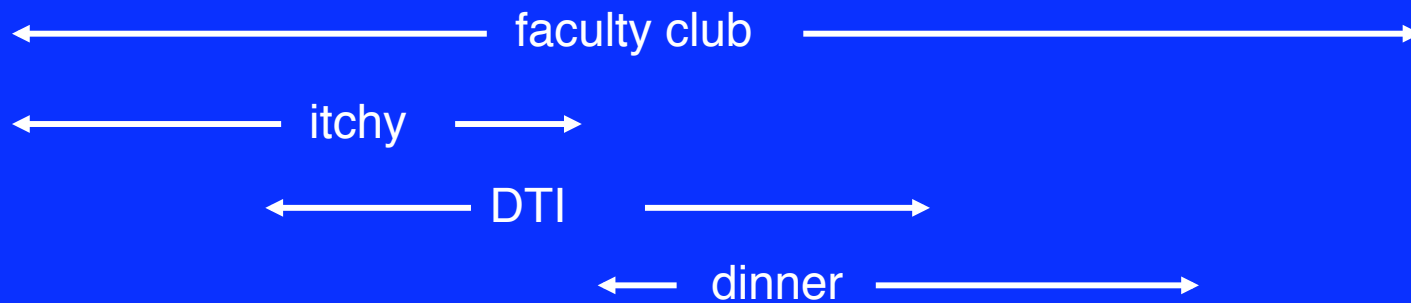
...delay...



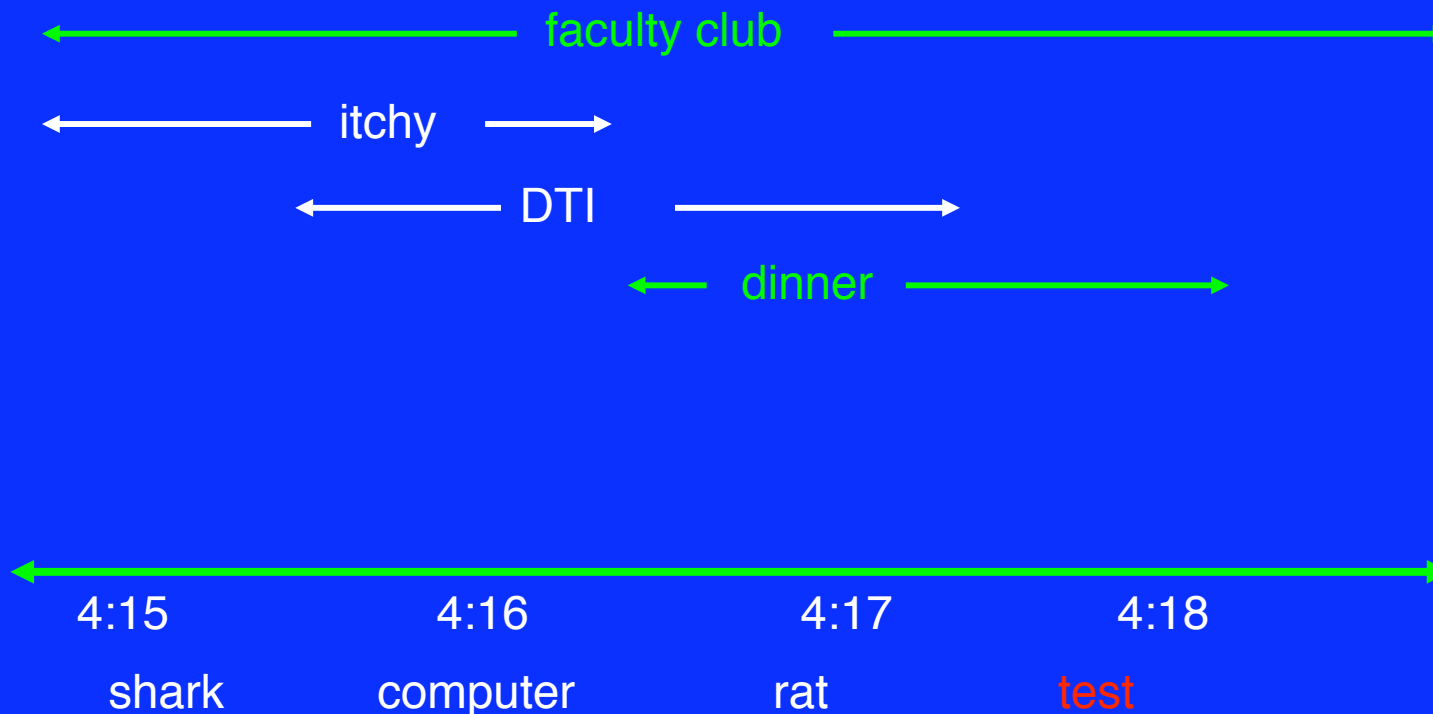
Drifting Mental Context



Drifting Mental Context



Context and Recall

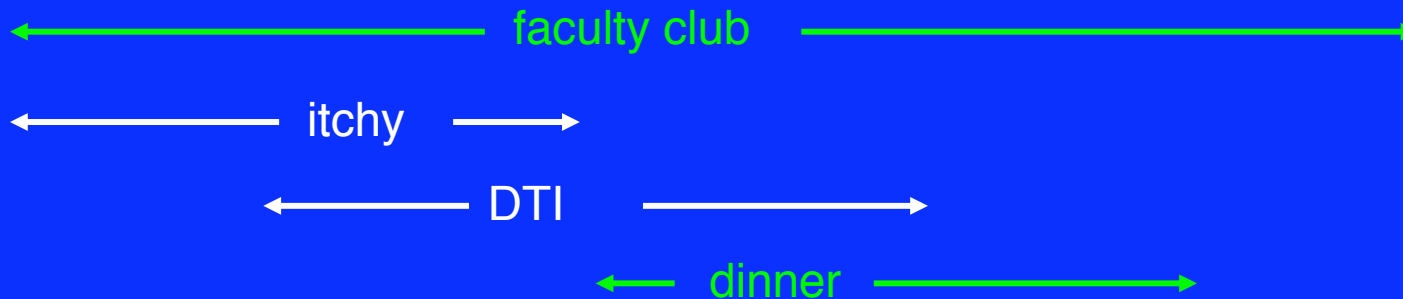


To recall the most recent list, cue with the current context:
“dinner, faculty club”

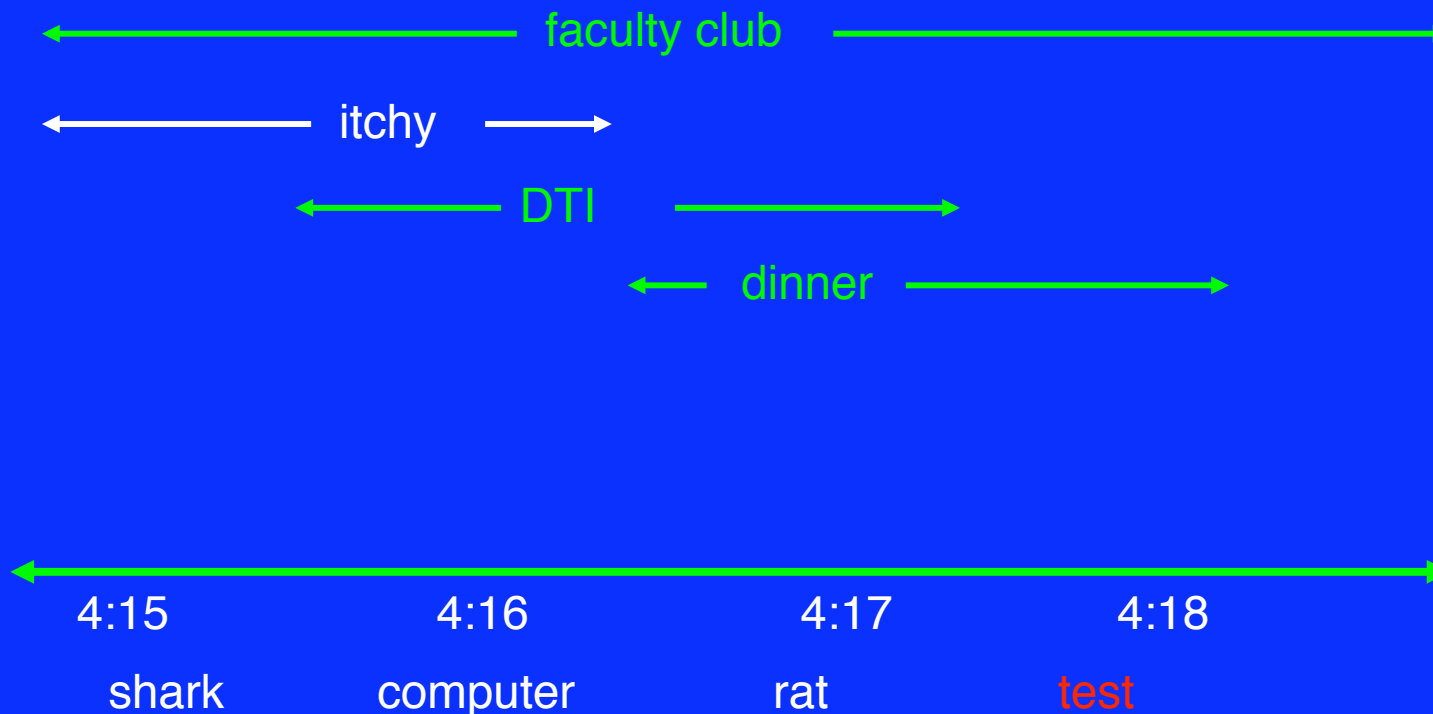
Given this cue, you end up recalling “rat”

You also recall **other contextual elements** associated with rat:
“DTI”

Context and Recall



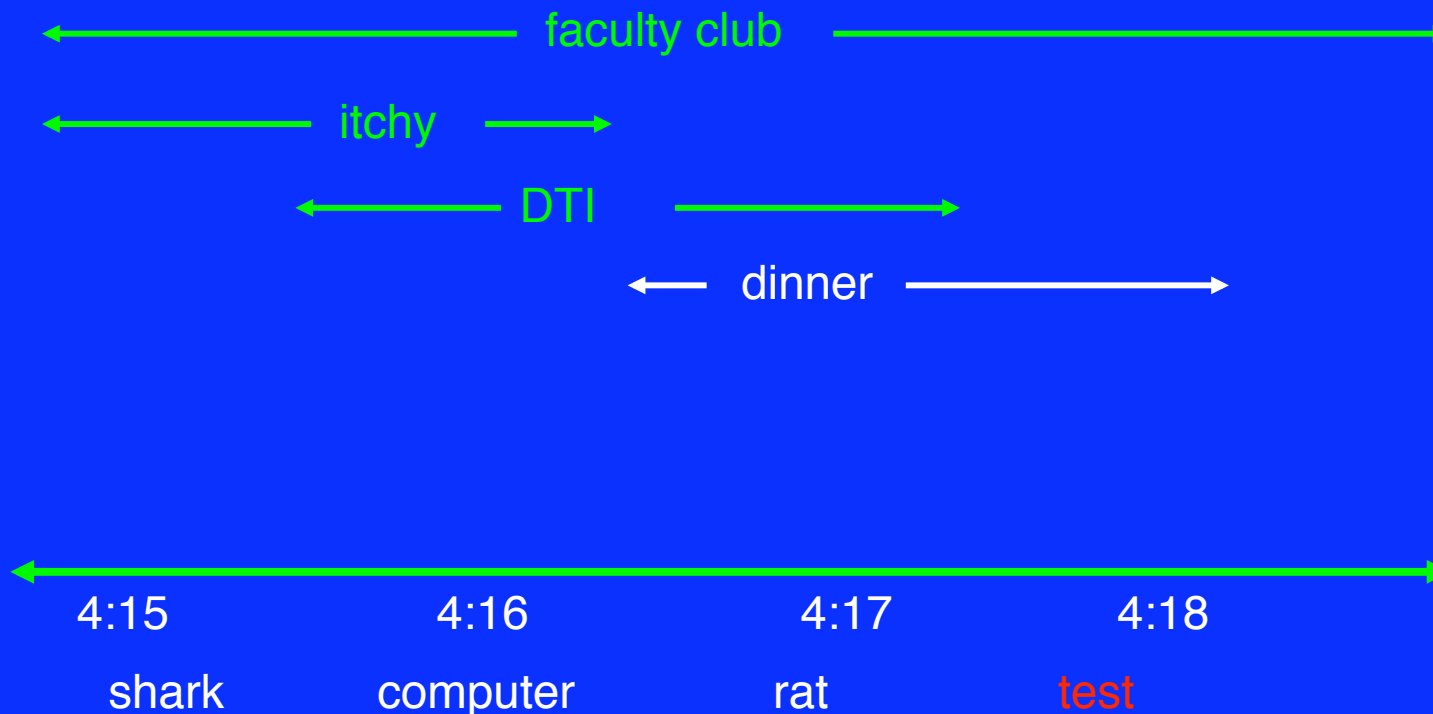
Context and Recall



Step 2: Take retrieved contextual elements and **incorporate them into your retrieval cue**: “dinner, DTI, faculty club”

With “DTI” in your retrieval cue, you can now recall “computer”, plus a new contextual element “itchy”

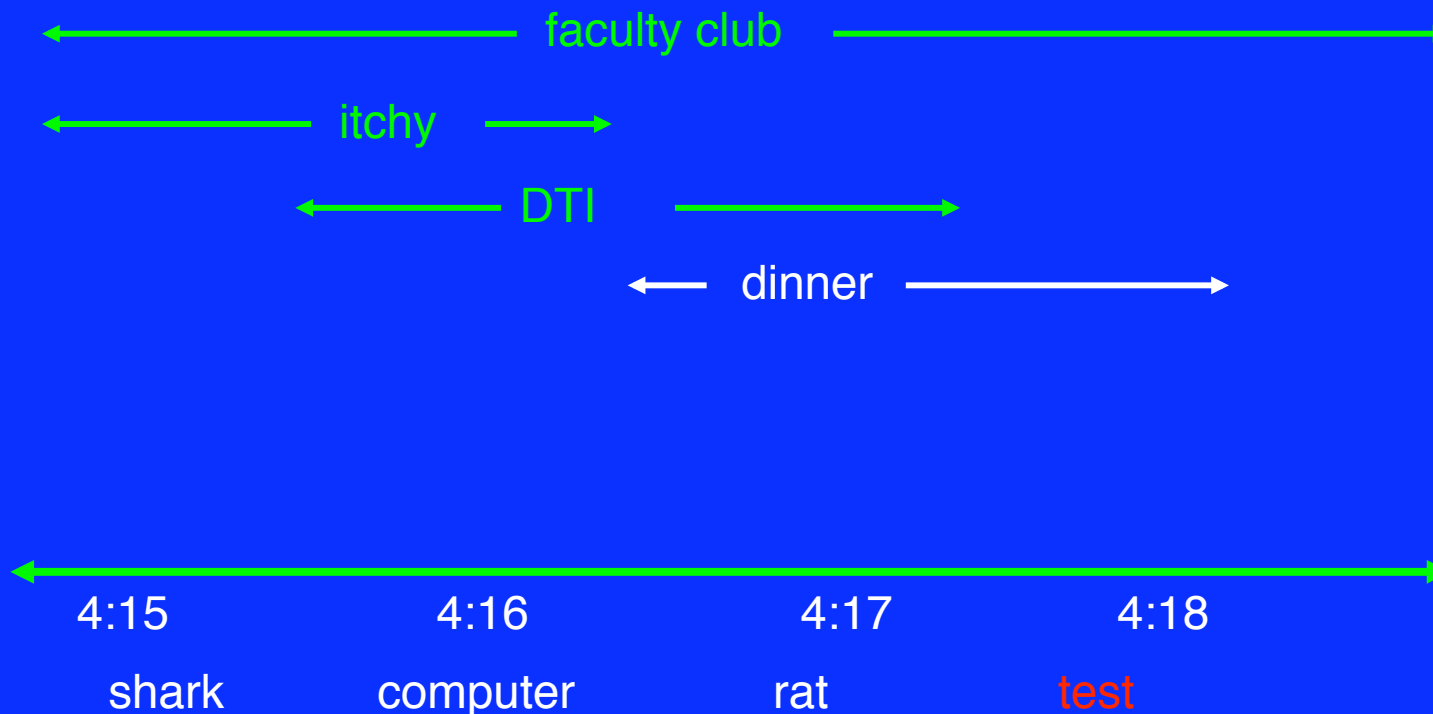
Context and Recall



Step 3: Incorporate “itchy” in your retrieval cue
Now you can recall “shark”

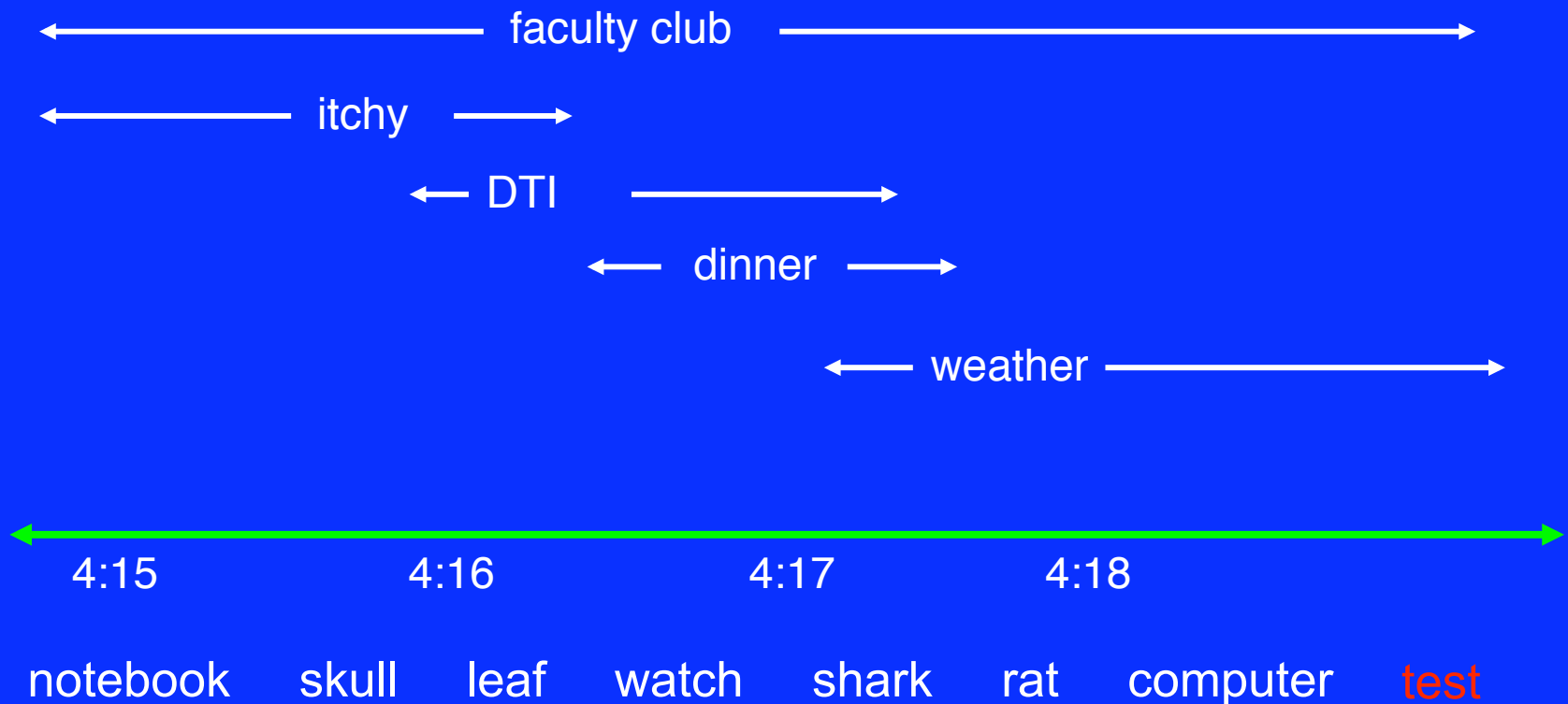
Using **retrieved context** as a retrieval cue allows you to bootstrap your way backwards in time...

Context and Recall



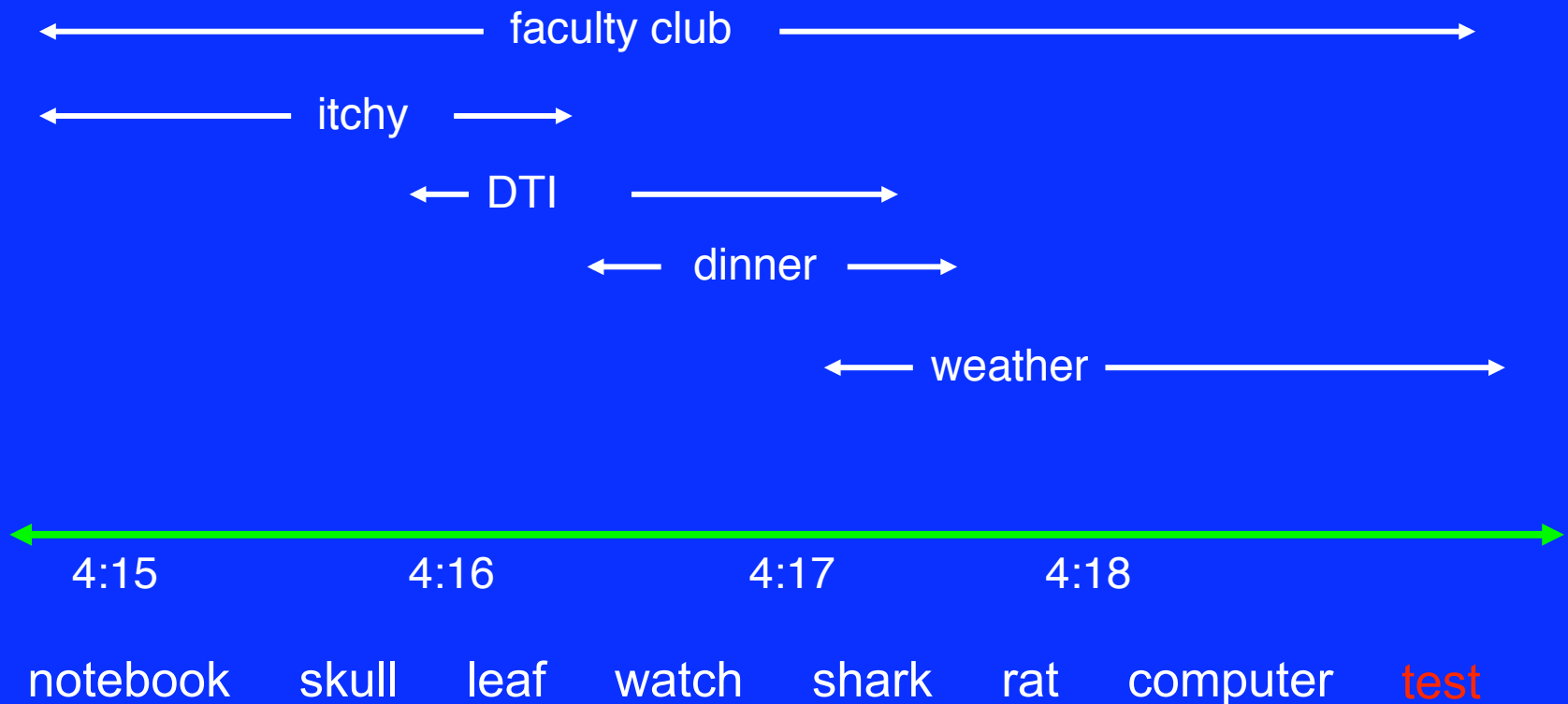
Key idea: Memory retrieval success depends on **contextual reinstatement**

Recency Judgments



This framework also suggests how people could make **judgments of recency**

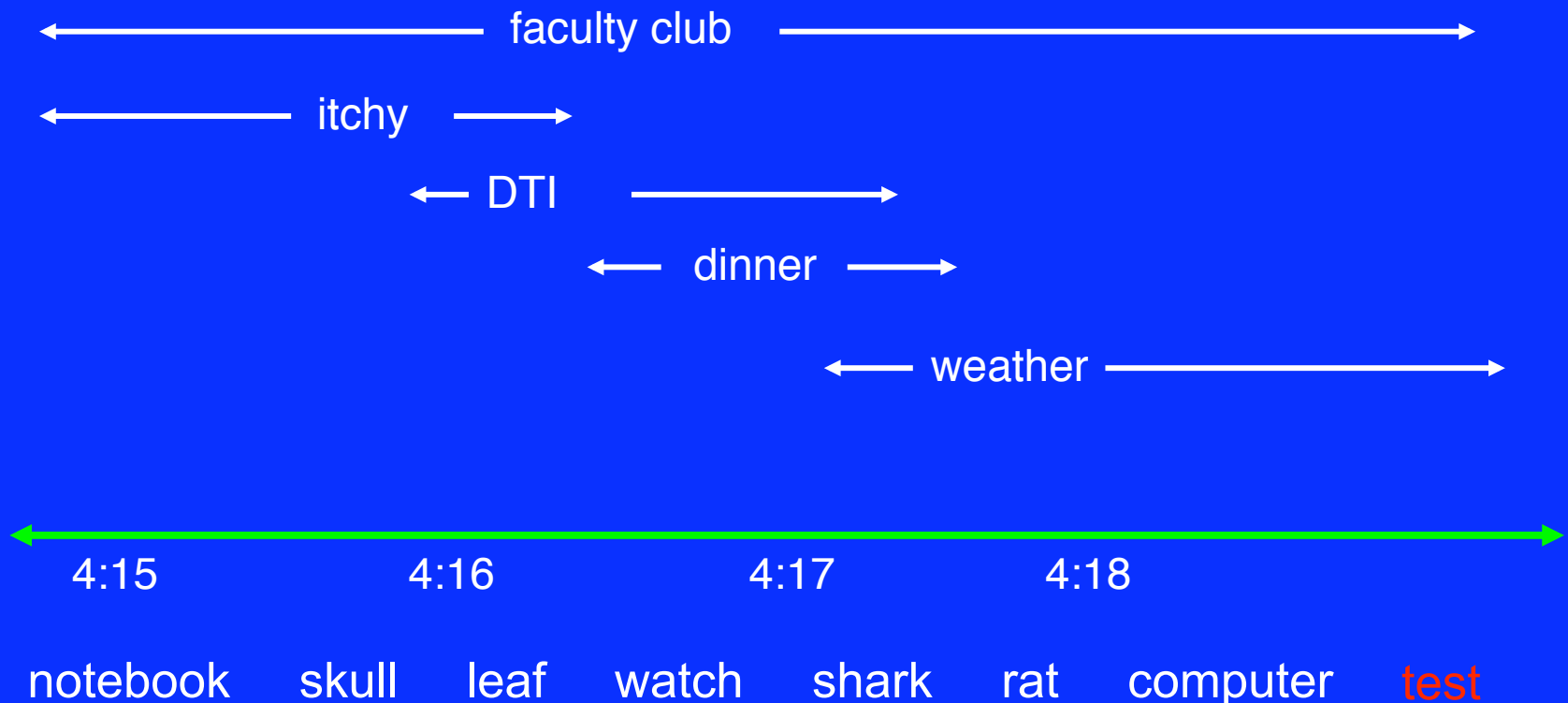
Recency Judgments



Which was presented more recently: shark or skull?

Use the words to cue for contextual info

Recency Judgments



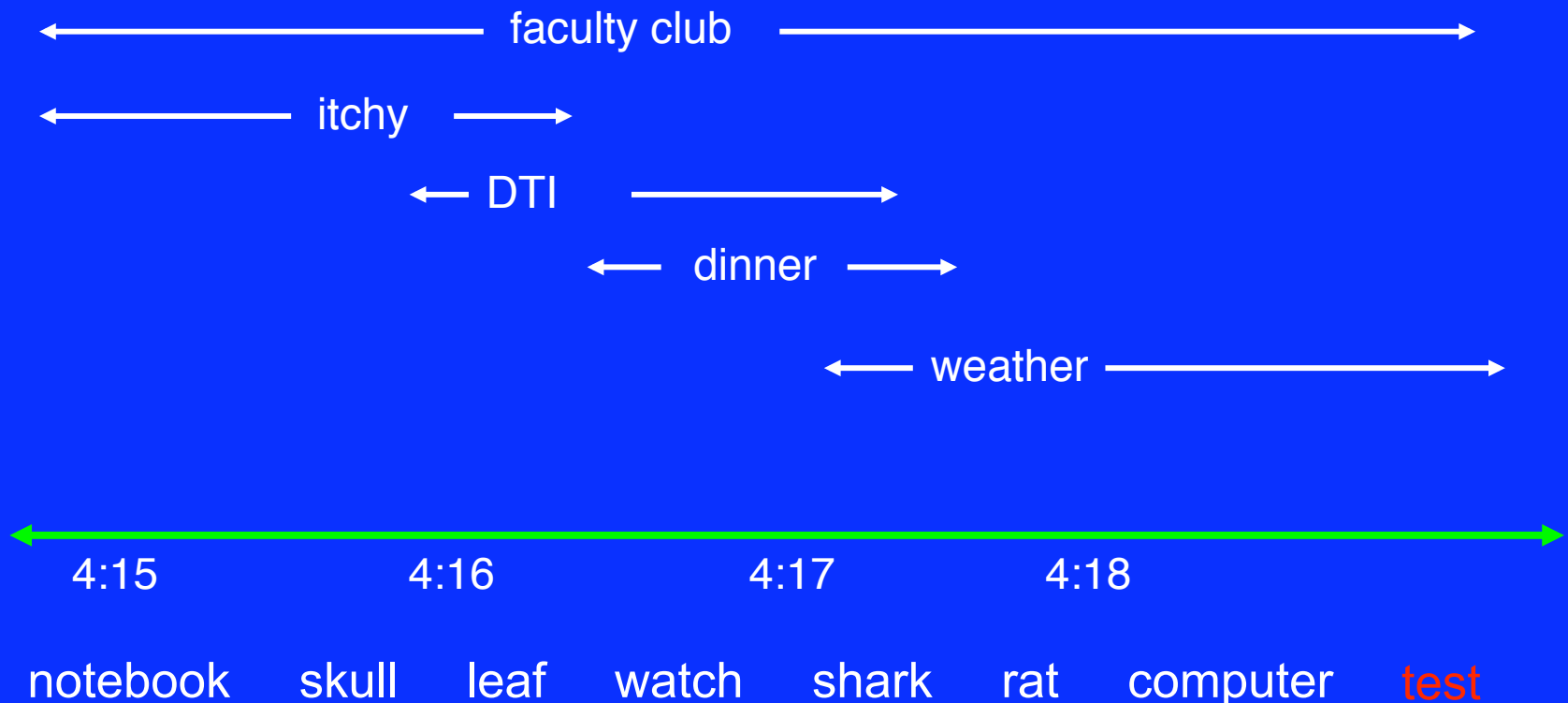
shark retrieves “dinner, weather, DTI, faculty club”

skull retrieves “DTI, itchy, faculty club”

compare retrieved context to current context:

“weather, faculty club”

Recency Judgments



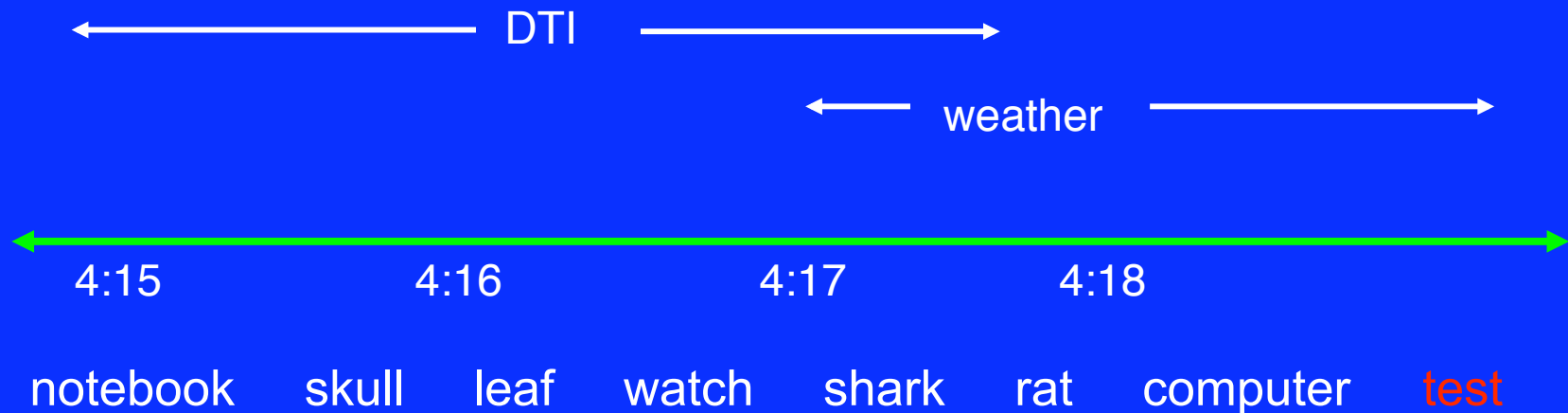
shark: “dinner, **weather**, DTI, **faculty club**”

skull: “DTI, itchy, **faculty club**”

current: “**weather**, **faculty club**”

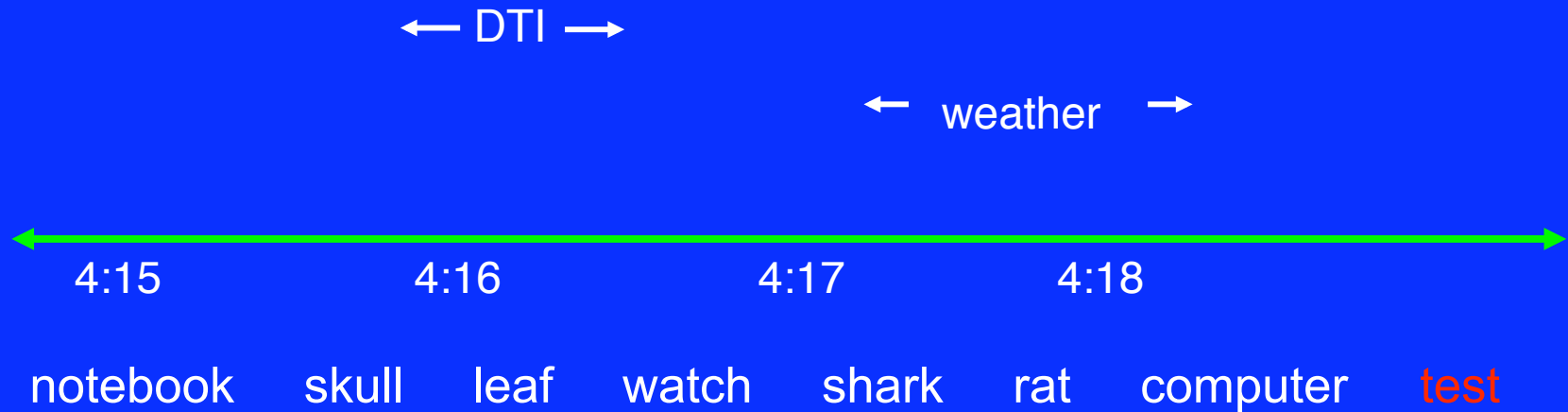
The shark context is **more similar** to the current context, so shark probably occurred more recently

Temporal Extent



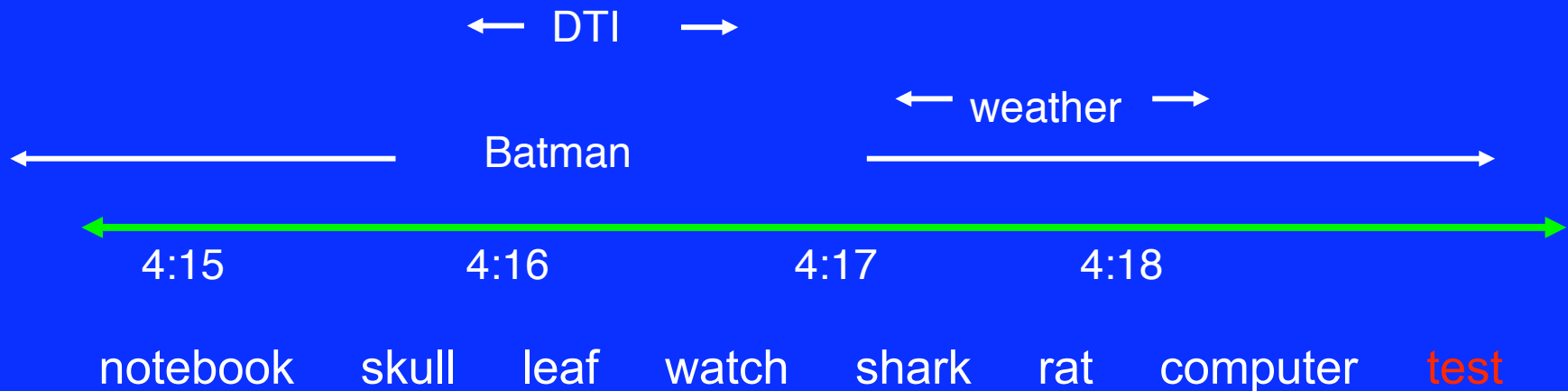
- Contextual representations are useful in cuing memory because of their **temporal extent**
- e.g., “weather” is useful in cuing “shark” at test because it extends temporally to cover both “shark” and “test”

Temporal Extent



- If contextual threads are too short, they aren't useful as memory cues

Temporal Extent

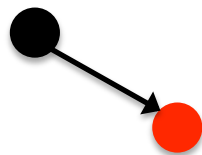


- If contextual threads persist for too long, they get **overloaded** and lose their efficacy as memory cues
- Memory retrieval in the brain is a **competitive** process
- Cues are effective if they **differentially** support some memories relative to others
- If a cue is linked (with uniform strength) to a large set of memories, it ceases to be an effective cue for the individual memories in the set

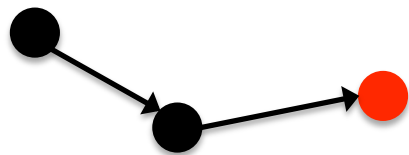
Context Space



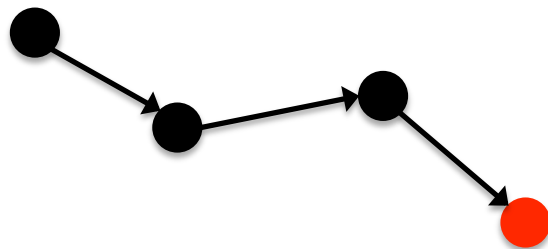
Context Space: Encoding



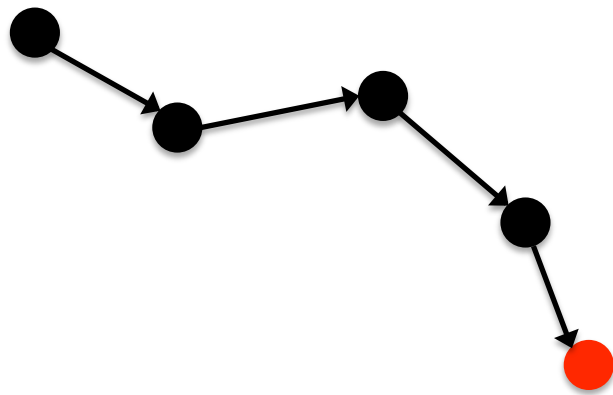
Context Space: Encoding



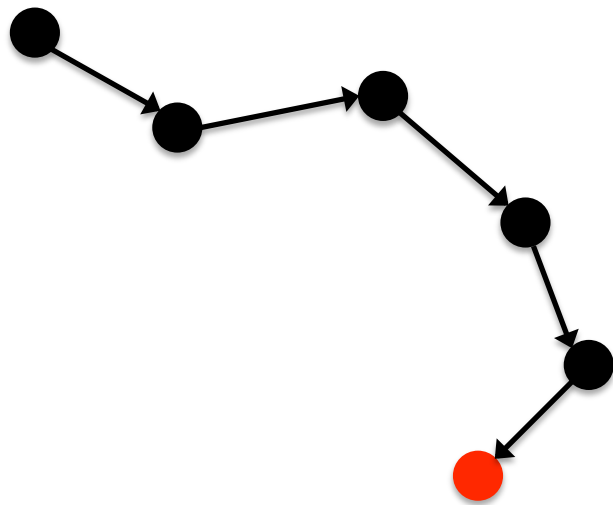
Context Space: Encoding



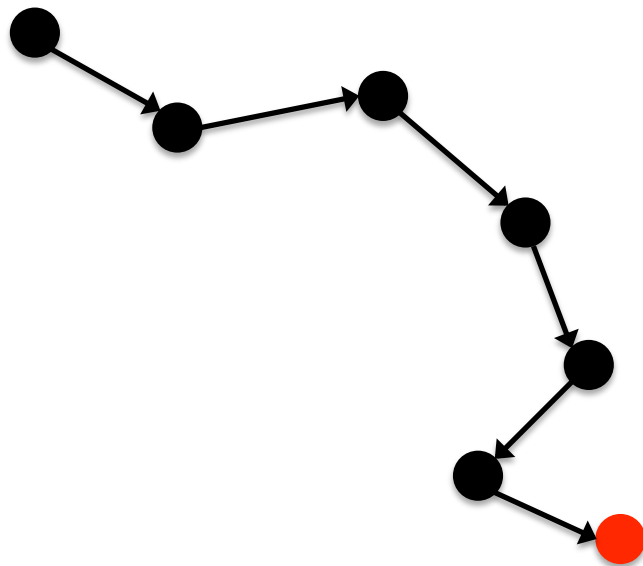
Context Space: Encoding



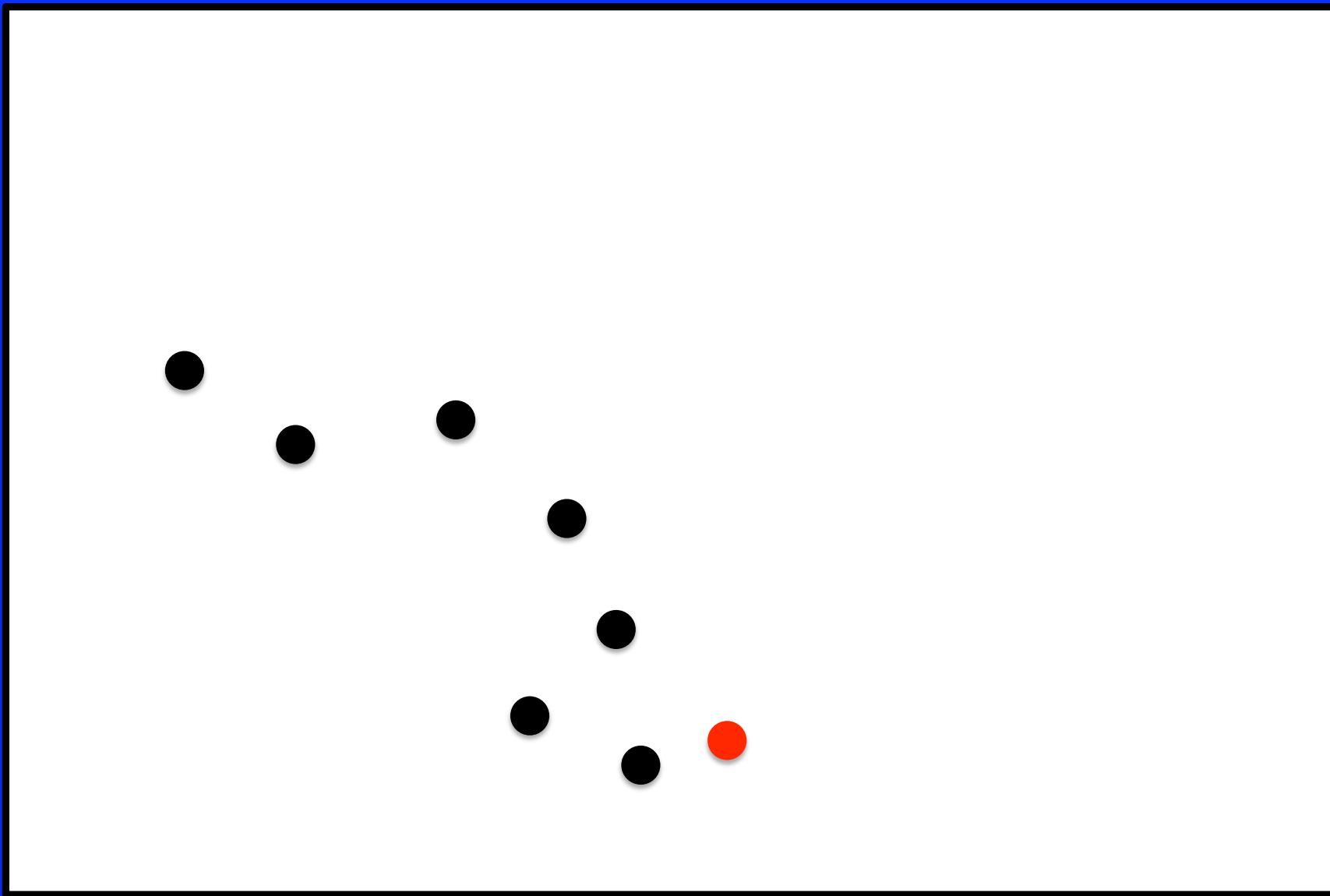
Context Space: Encoding



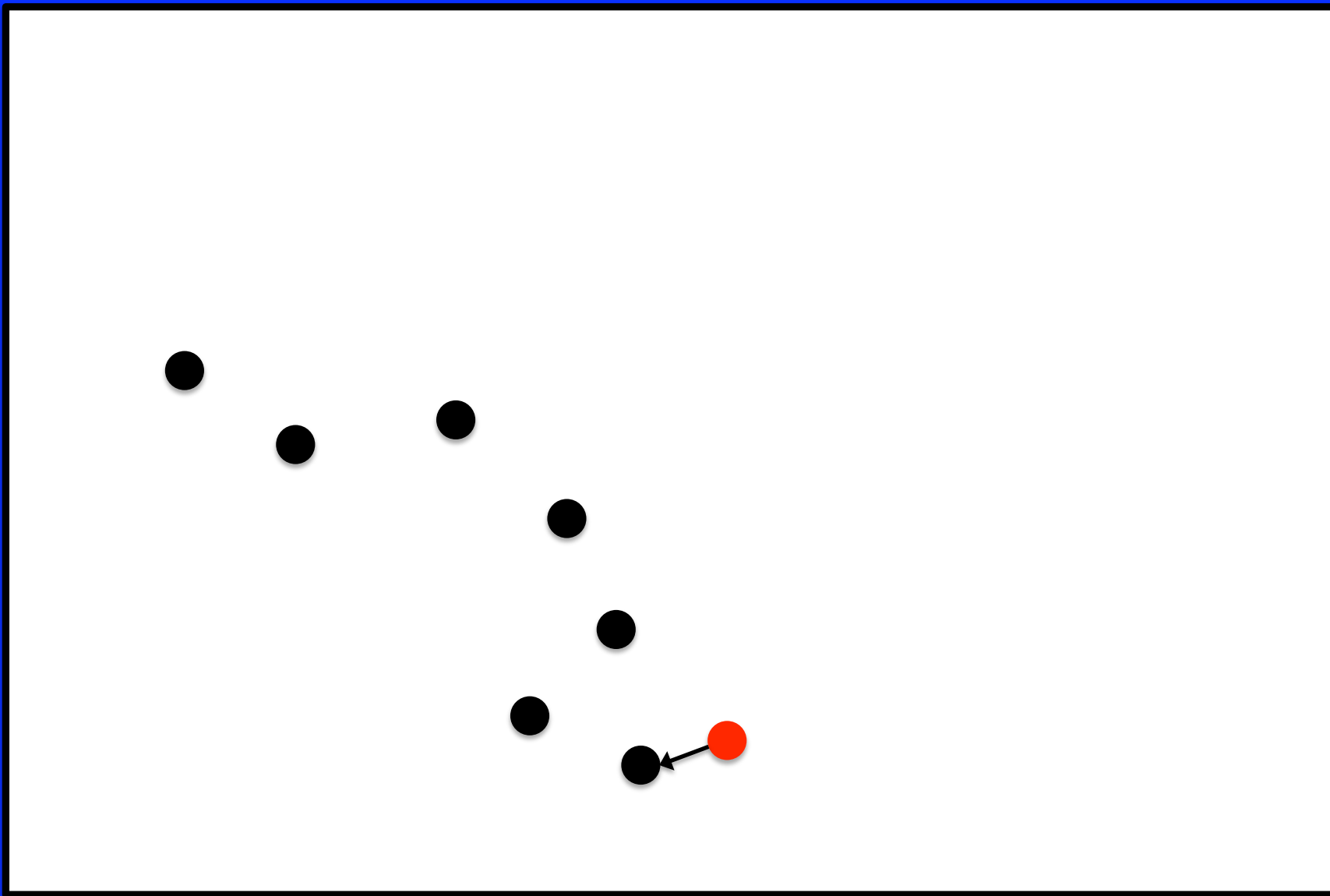
Context Space: Encoding



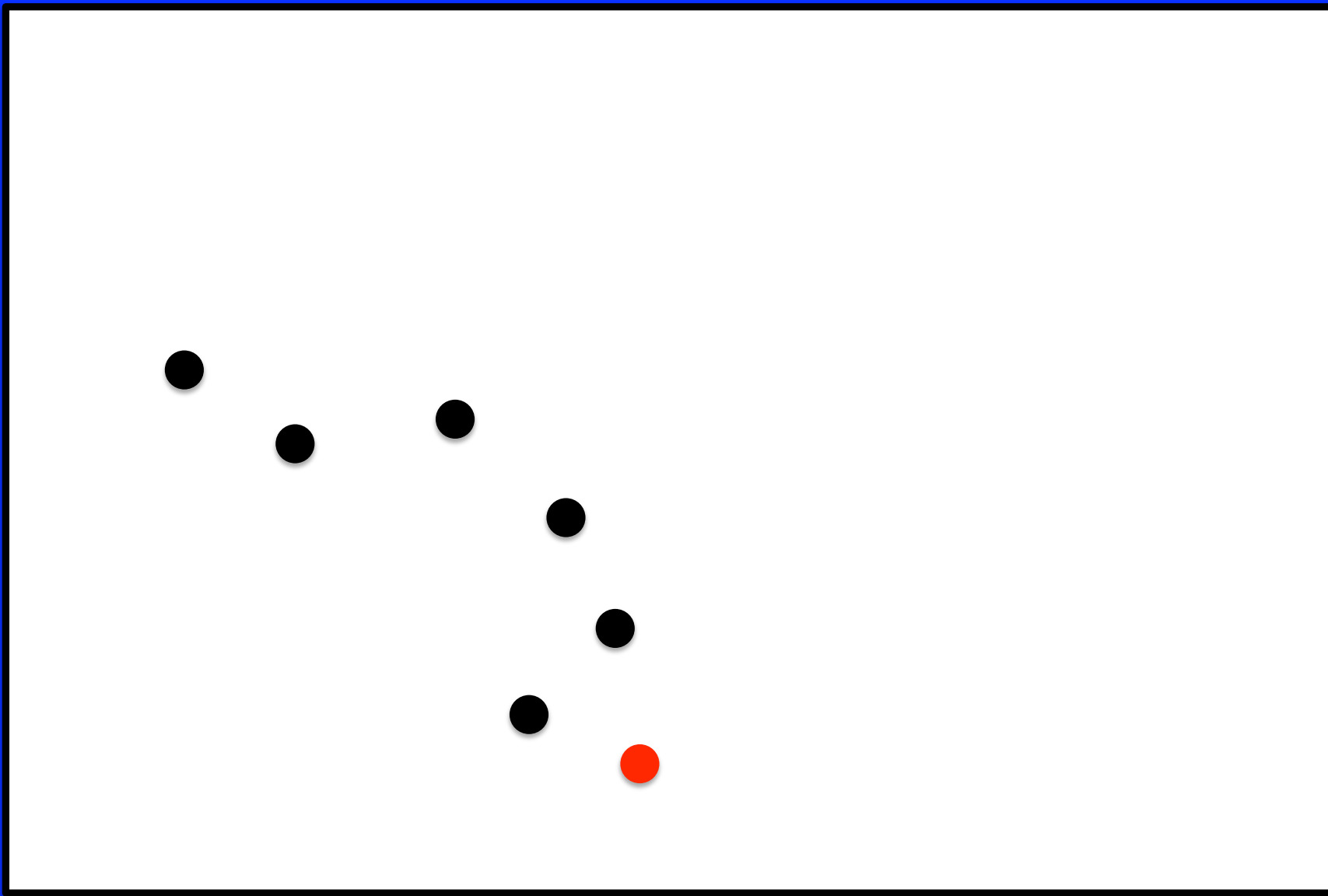
Context Space: Retrieval



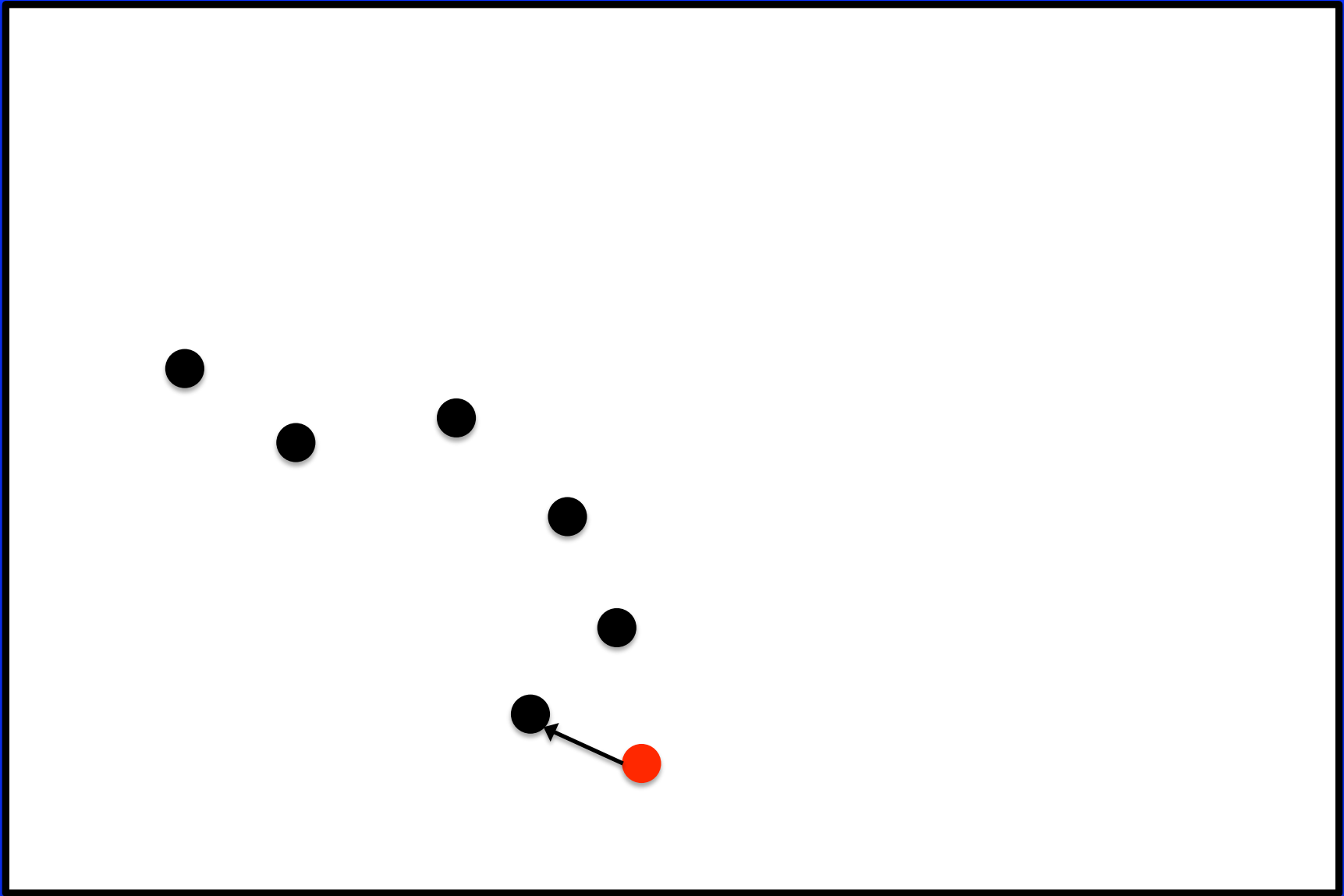
Context Space: Retrieval



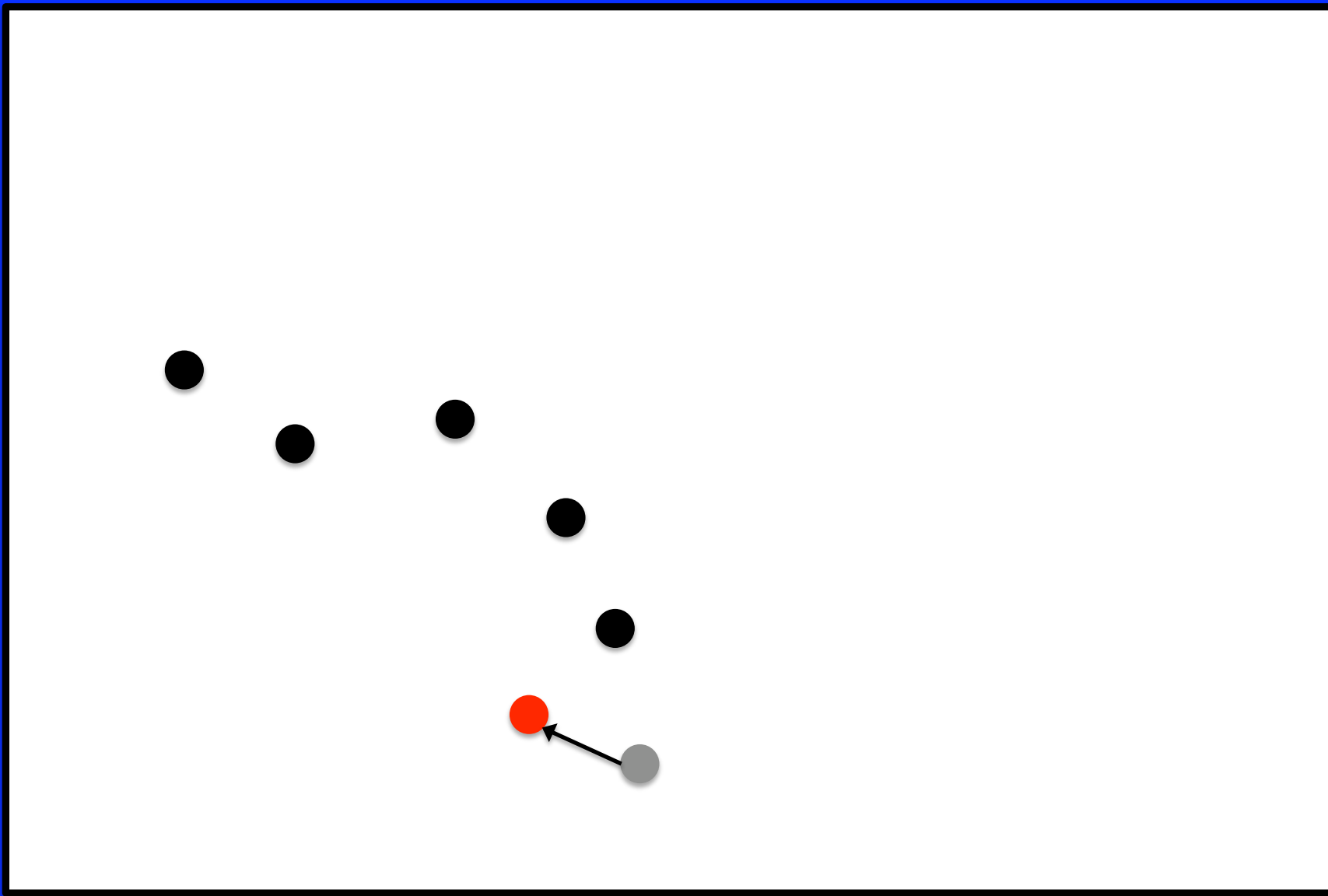
Context Space: Retrieval



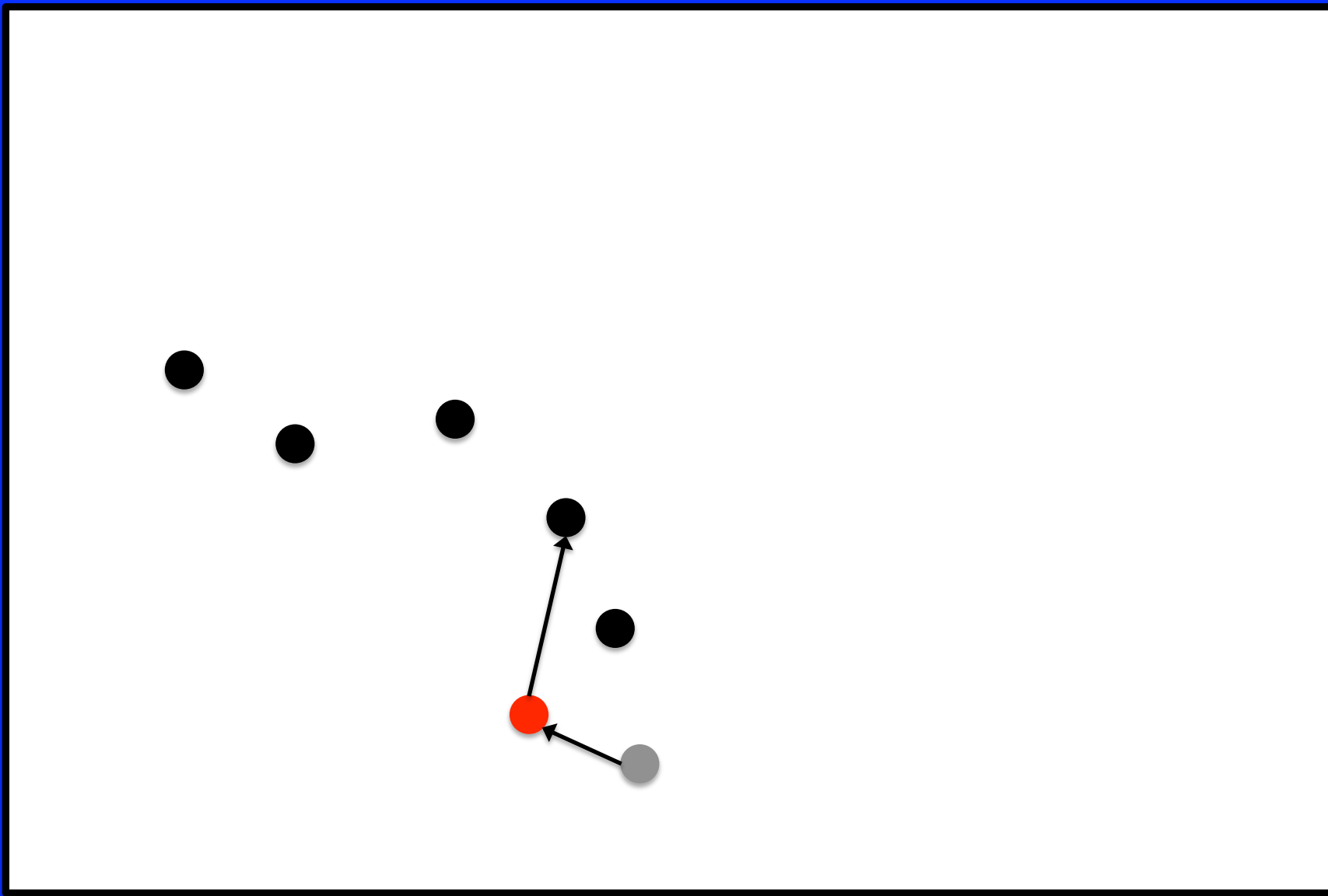
Context Space: Retrieval



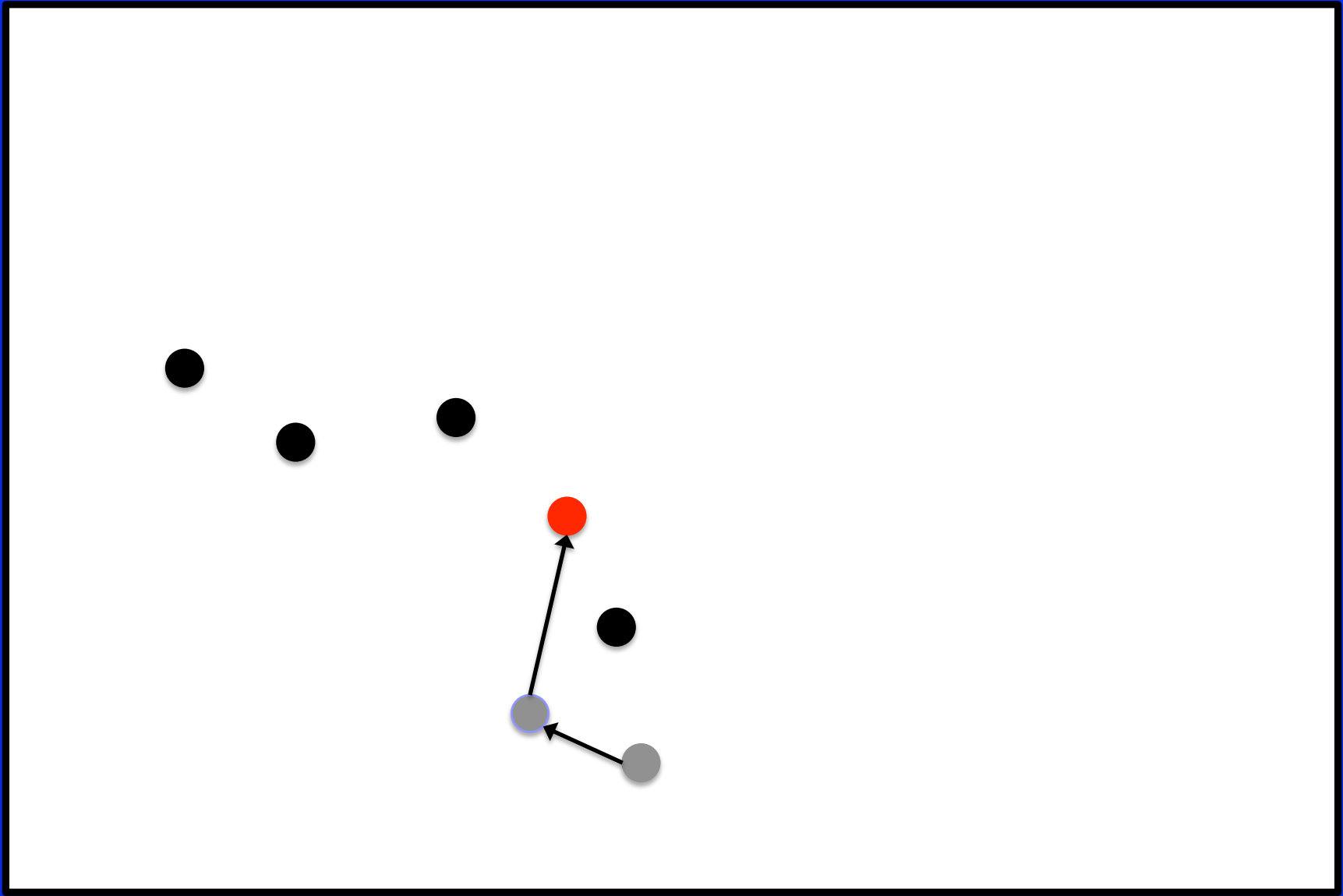
Context Space: Retrieval



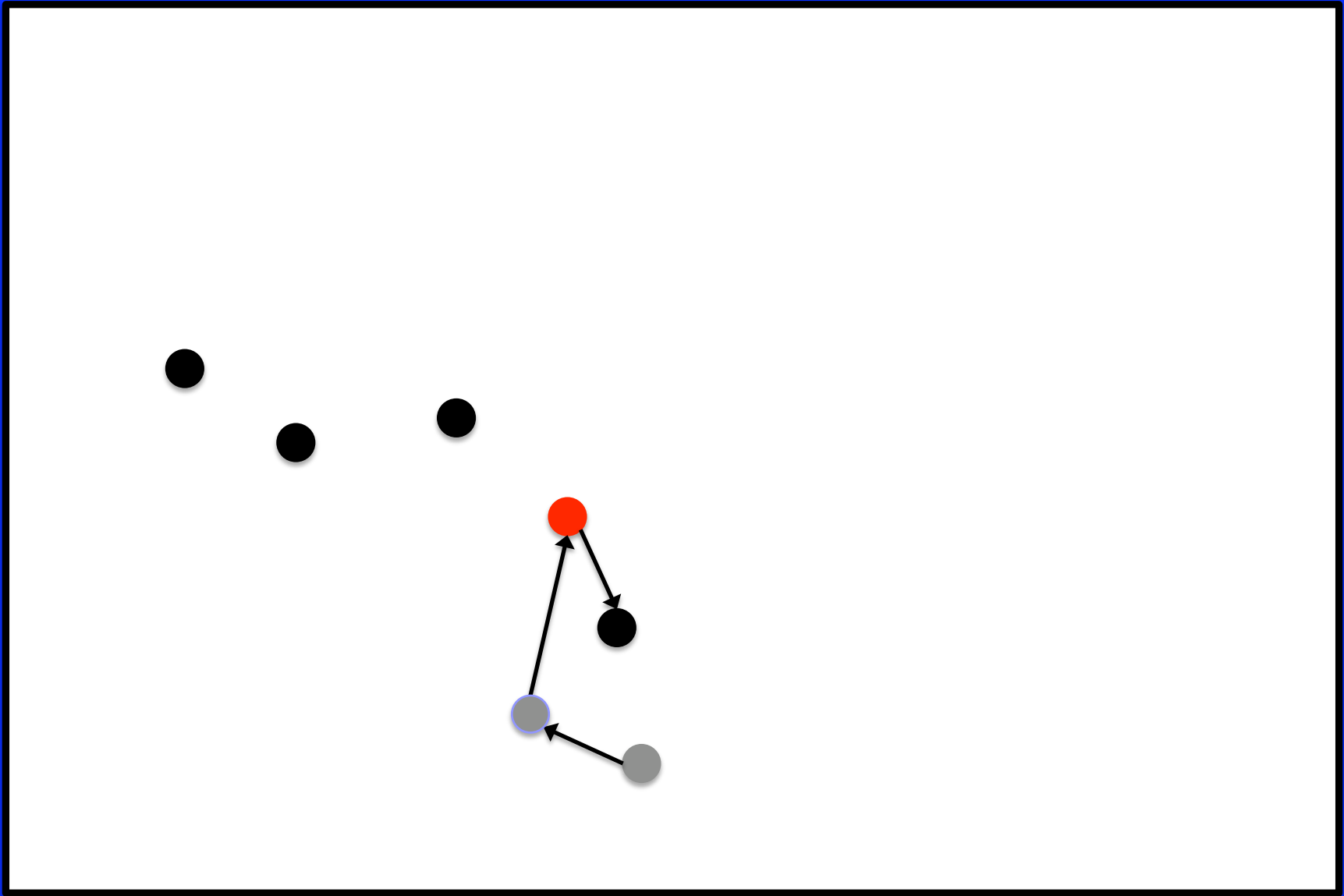
Context Space: Retrieval



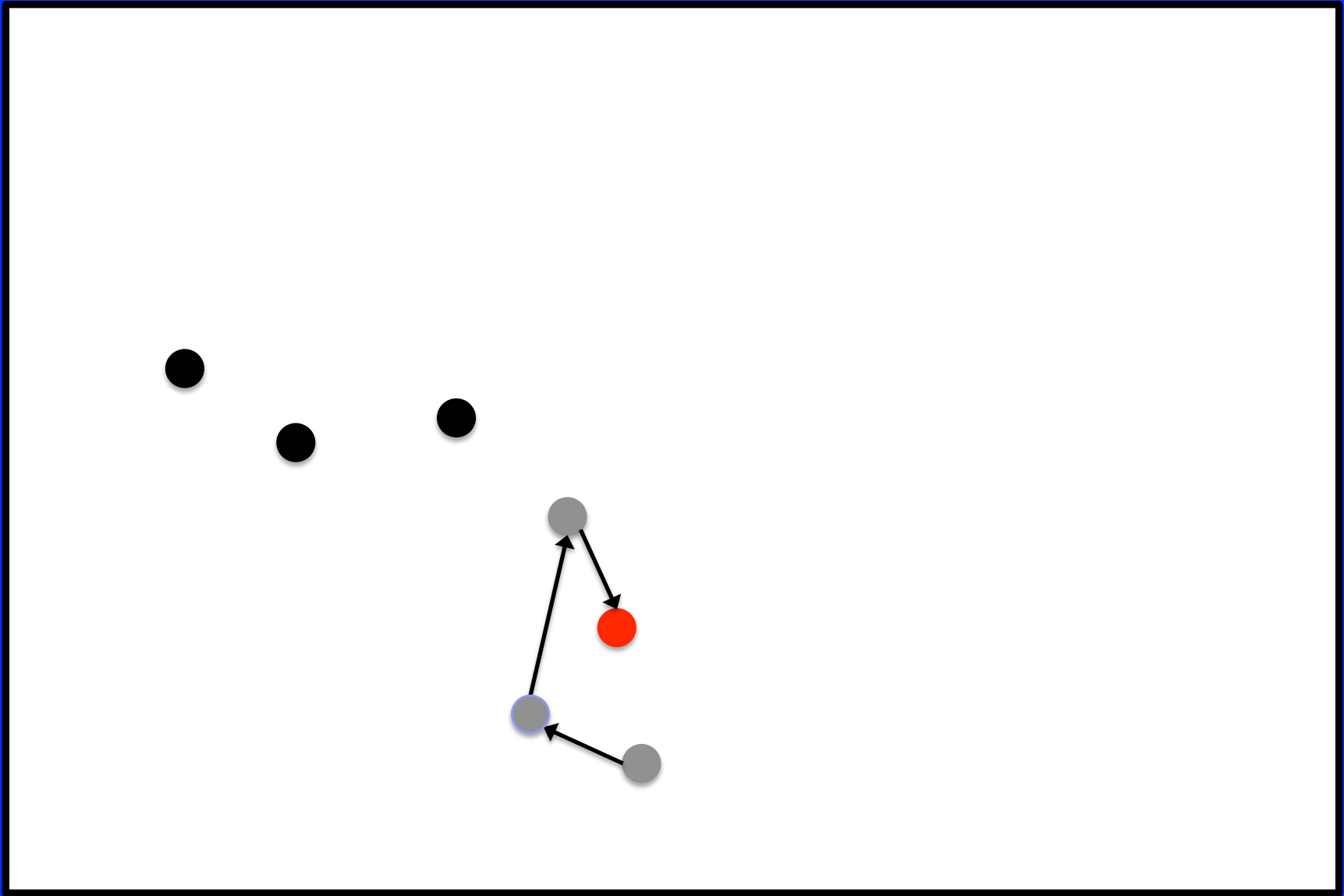
Context Space: Retrieval



Context Space: Retrieval



Context Space: Retrieval



Computational Models of Memory Search

- Psychologists have developed explicit computational models of memory search
 - the Temporal Context Model (TCM: Howard & Kahana, 2002; Sederberg et al., in press, *Psychological Review*)
 - the Context Maintenance and Retrieval model (CMR; Polyn, Norman, & Kahana, submitted)
- These computational models operationalize context as a slowly drifting vector. The context vector is associated with item vectors, such that items can trigger contextual retrieval and vice-versa.

Computational Models of Memory Search

- The models generate extremely detailed predictions about the trajectory of the context vector at encoding and retrieval, and the effects of contextual drift on behavior
- We can test some of these predictions by looking at behavioral data, but this is very indirect...
 - If subjects don't behave as predicted, it's difficult to know what went wrong
 - We don't know exactly how (or if) the context vector deviated from the predicted trajectory
- To properly evaluate these theories, we need to develop methods for **directly reading out** the state of the context vector based on brain data

fMRI studies

- Basic logic:
- Present items in different contexts at study
- Train a classifier (on study-phase data) to recognize the neural correlates of these contexts
- Measure reinstatement of these contexts at test

Tracking Memory Search (Polyn et al., 2005)

- Memory experiment: Subjects study of 3 types of stimuli

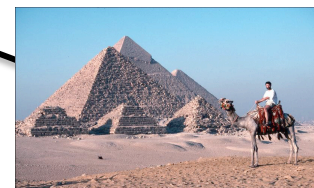


- Recall test: Recall items from all 3 categories, in any order
- Hypothesis: To recall a particular category, subjects try to reinstate the appropriate context from the study phase
- Concretely: To recall faces, subjects try to **make their brain state at test** resemble their brain state when they were studying faces
- If subjects succeed at recapturing their brain state from the study phase, this will trigger recall of specific studied items...

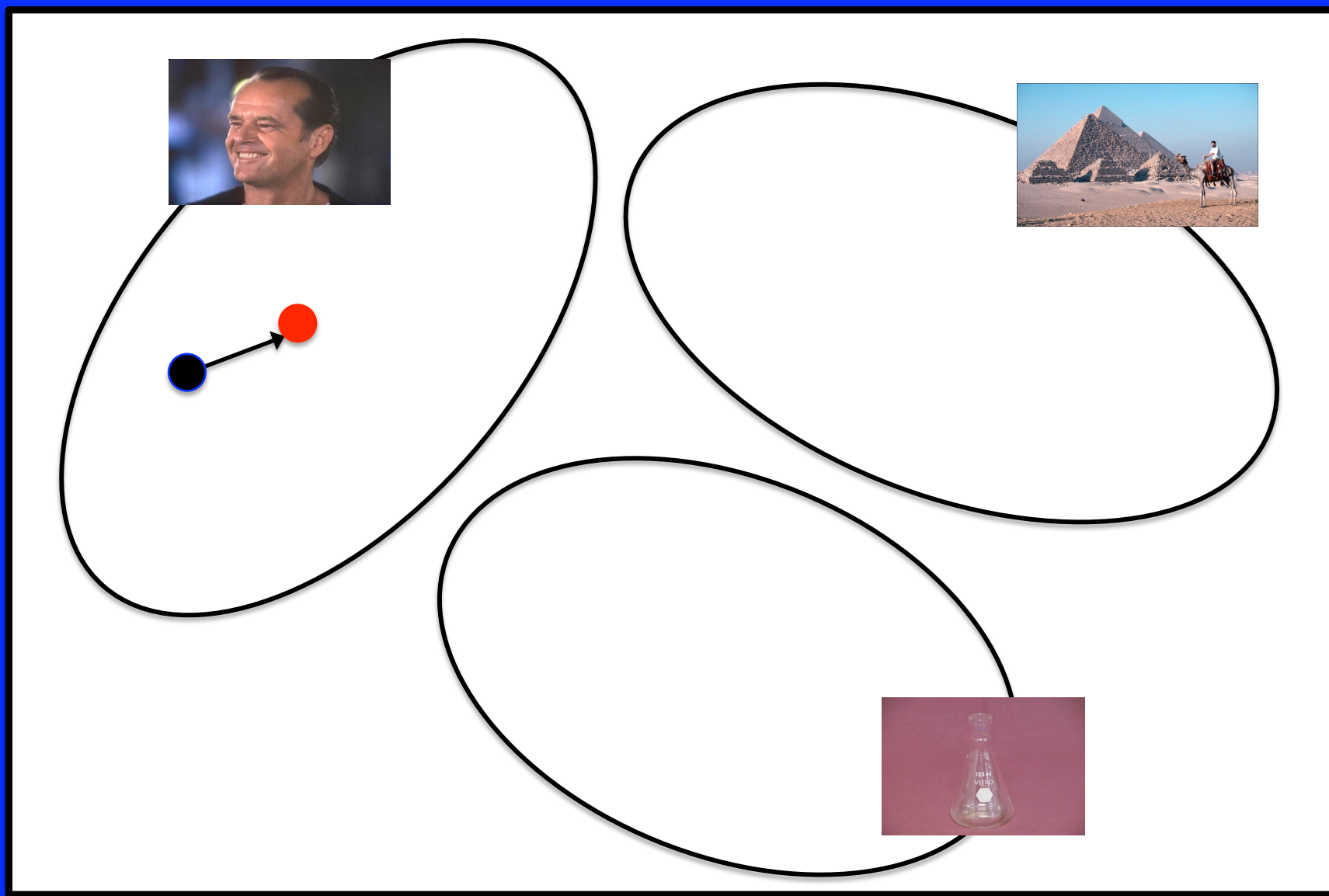
Context Space



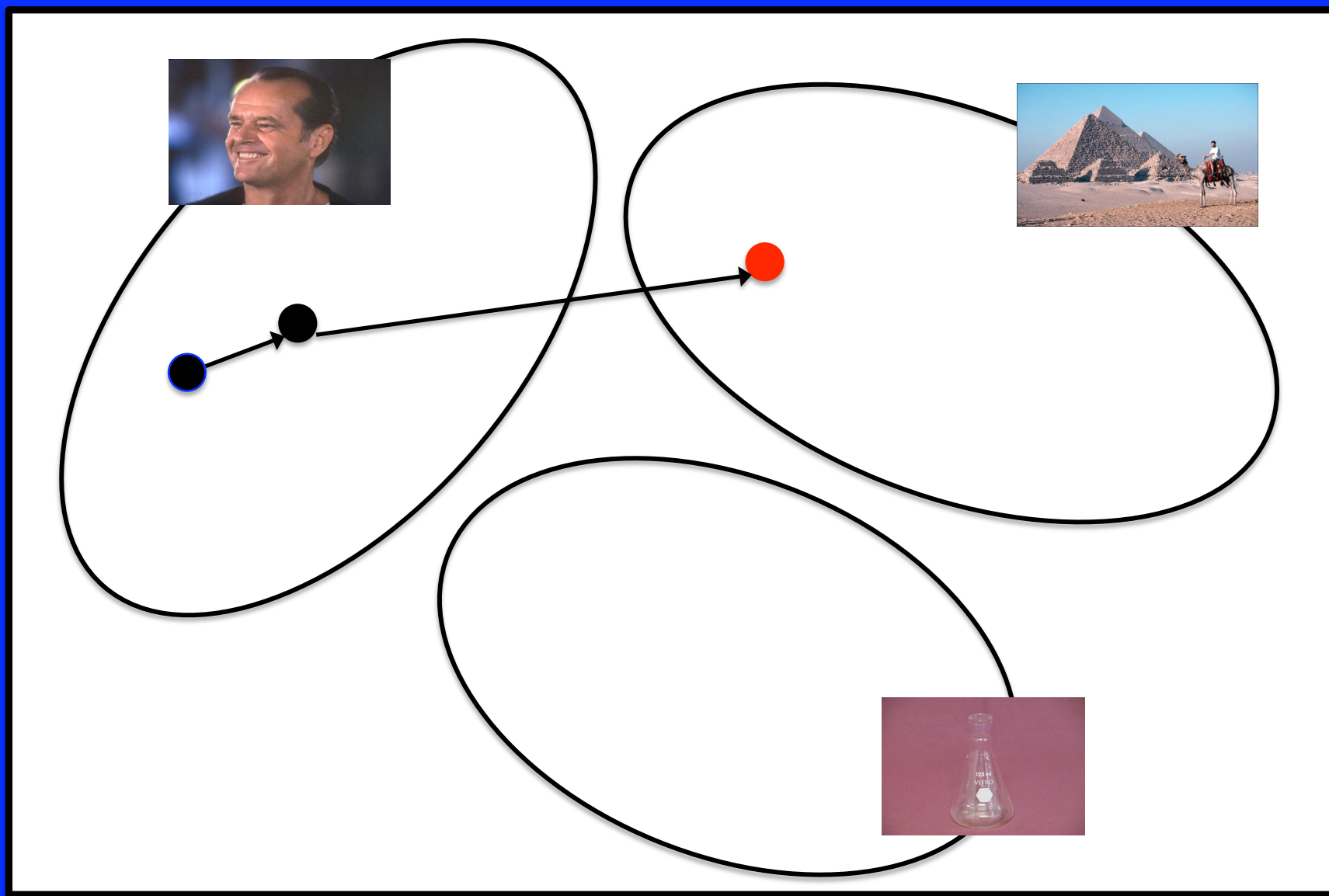
Context Space: Encoding



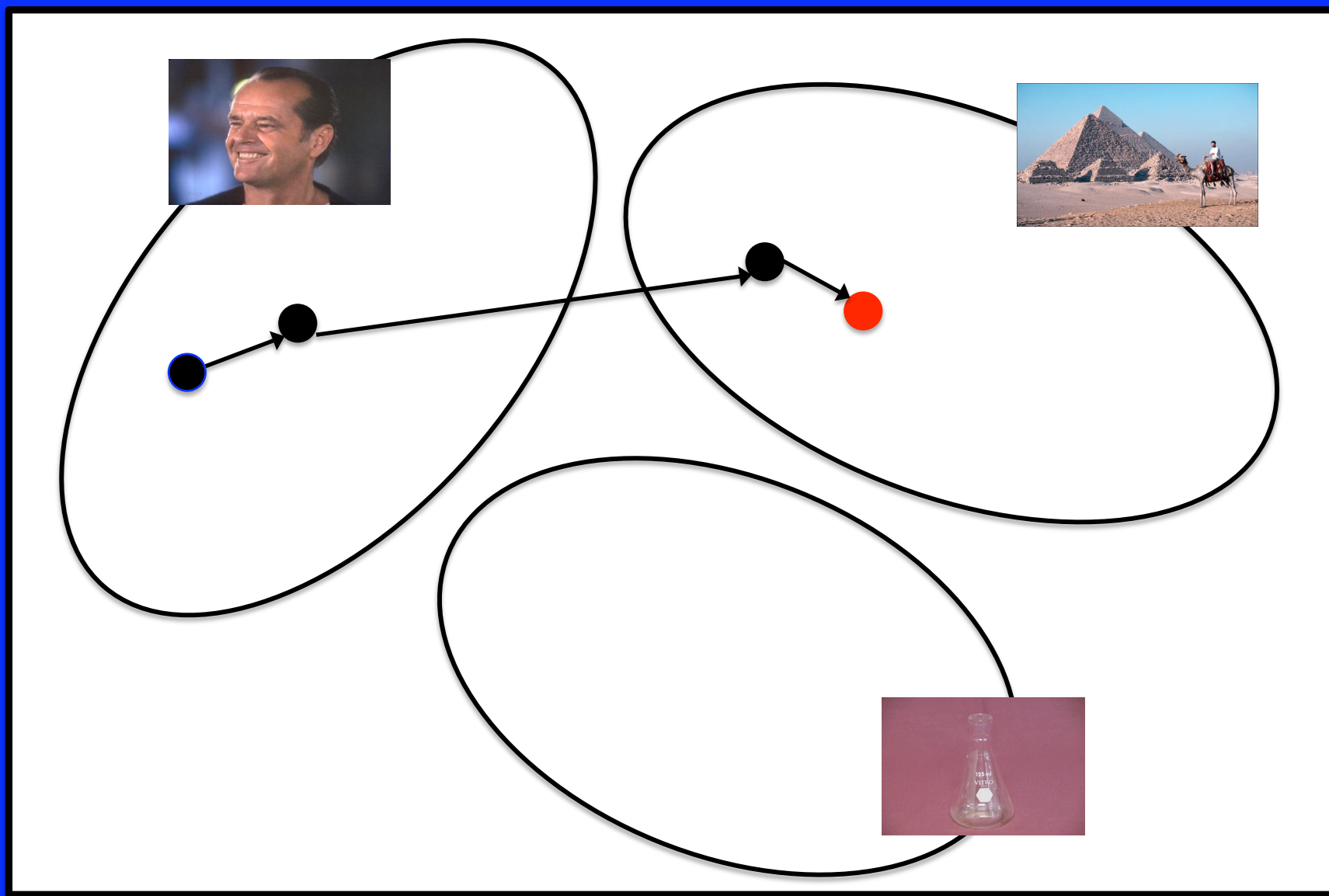
Context Space: Encoding



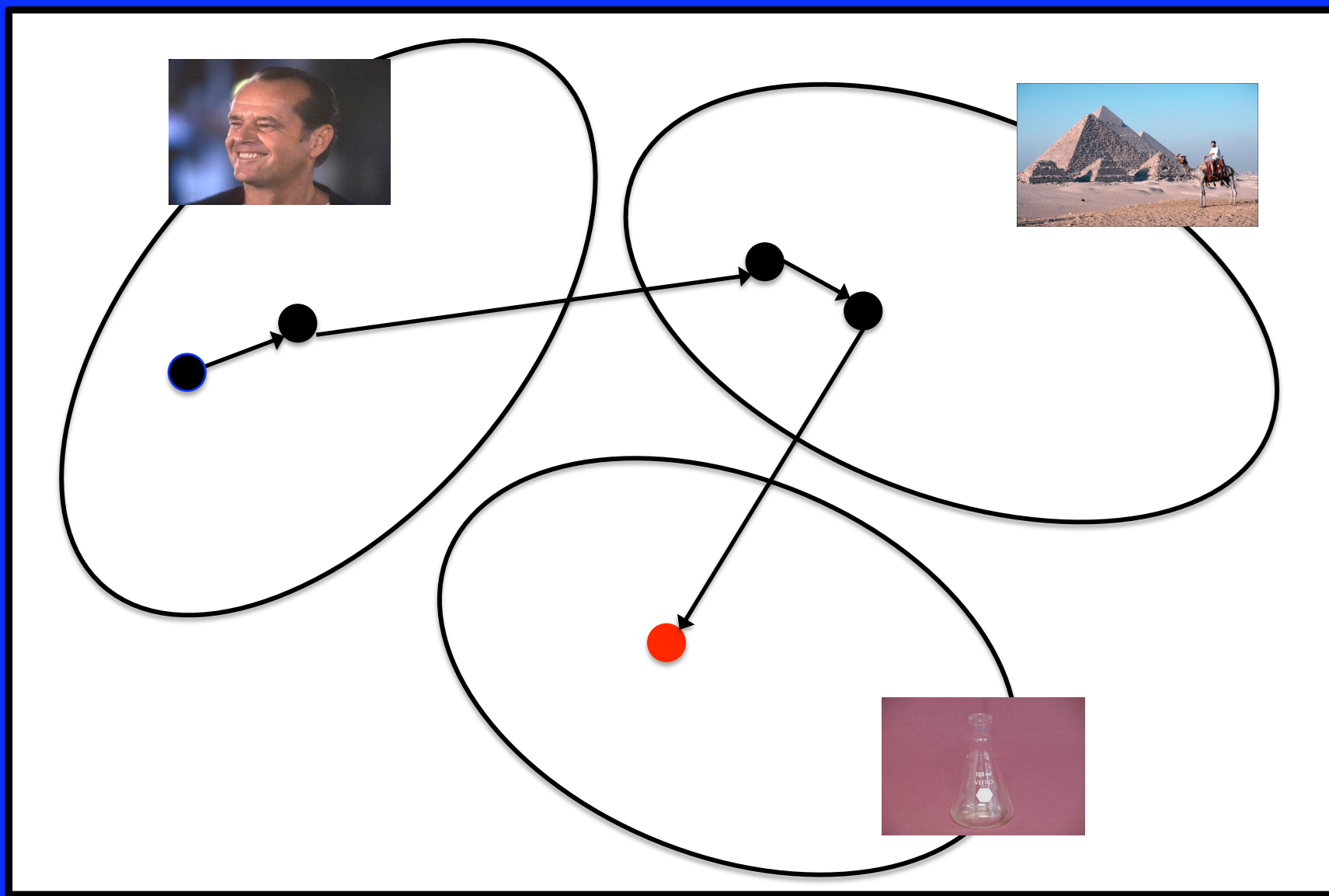
Context Space: Encoding



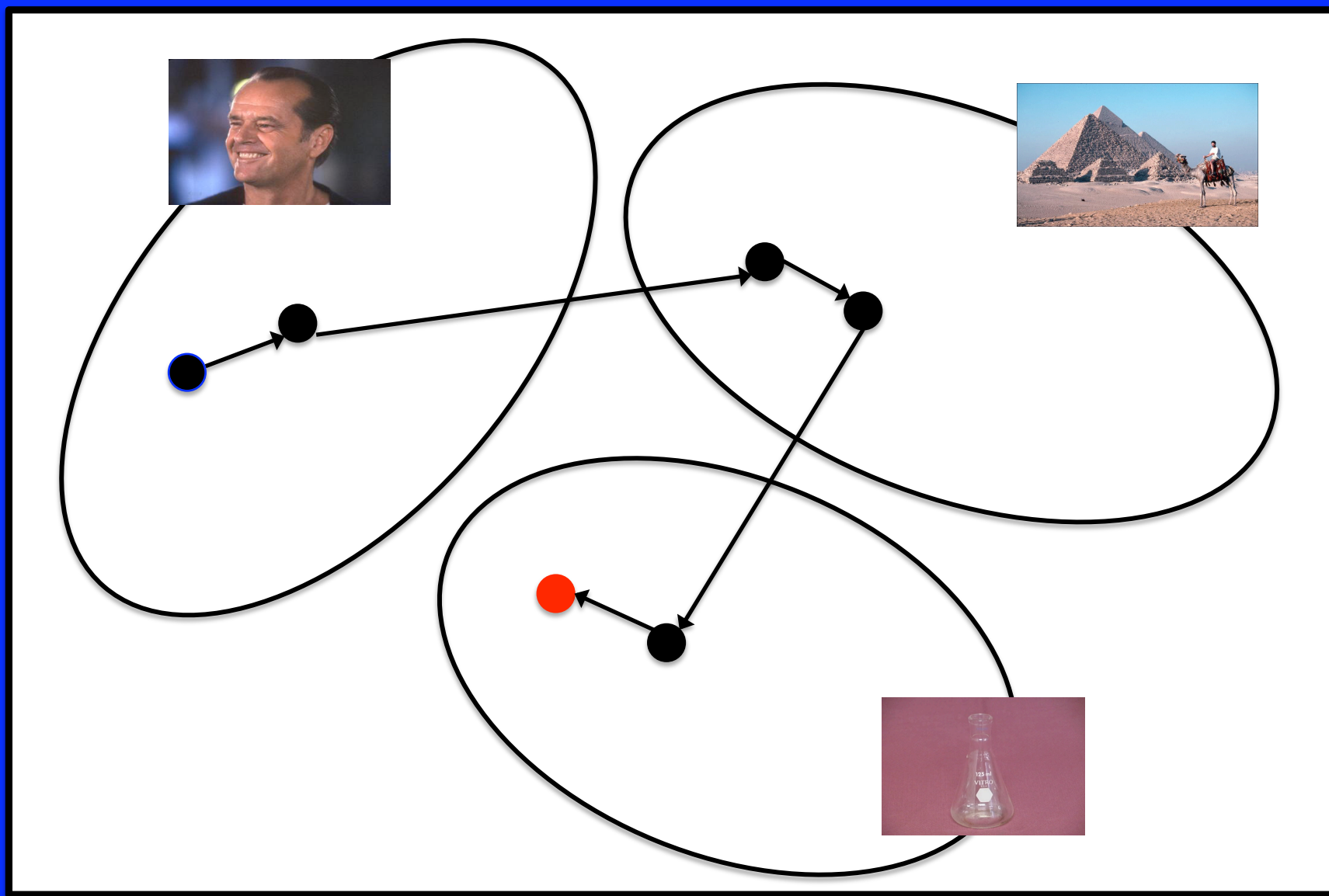
Context Space: Encoding



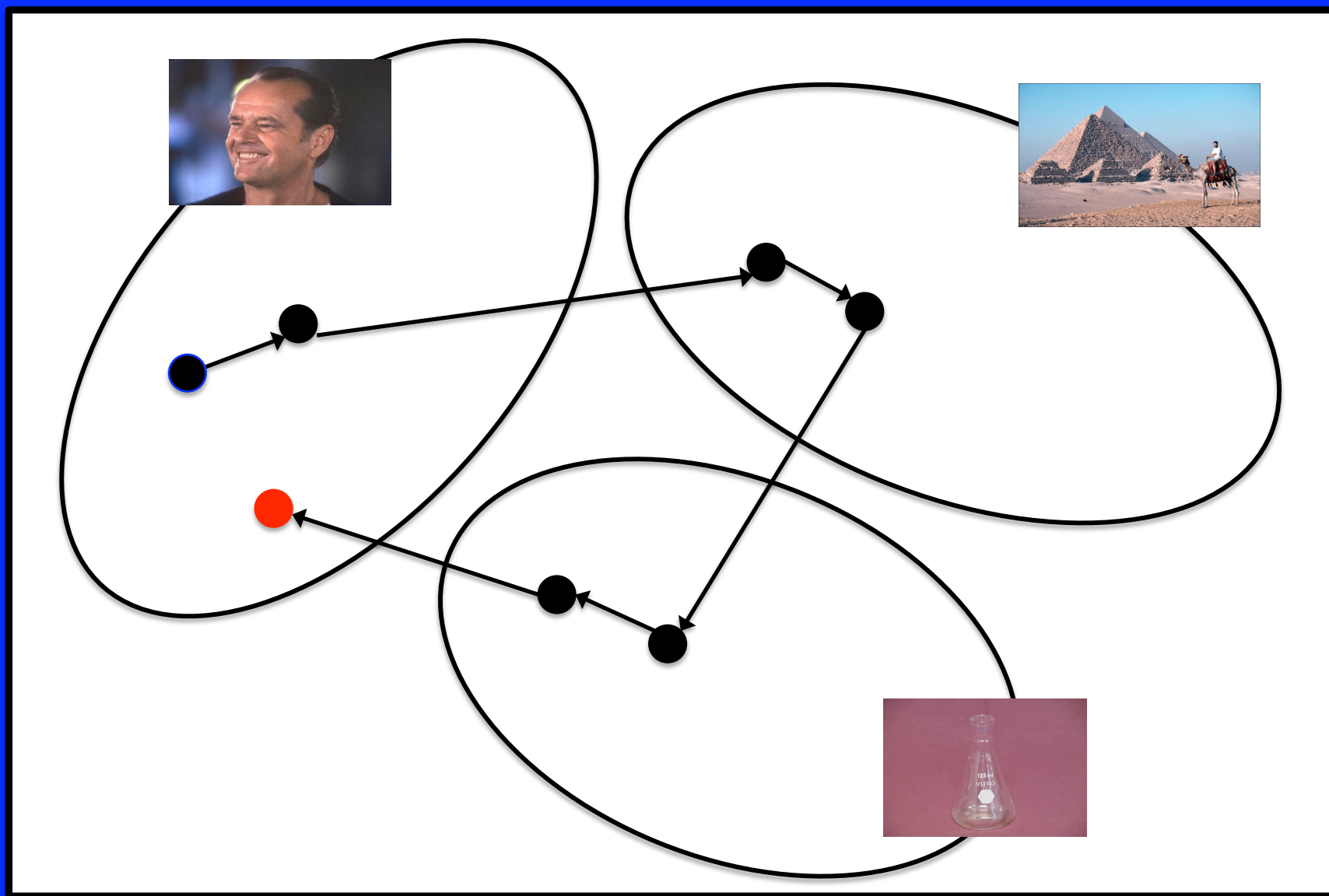
Context Space: Encoding



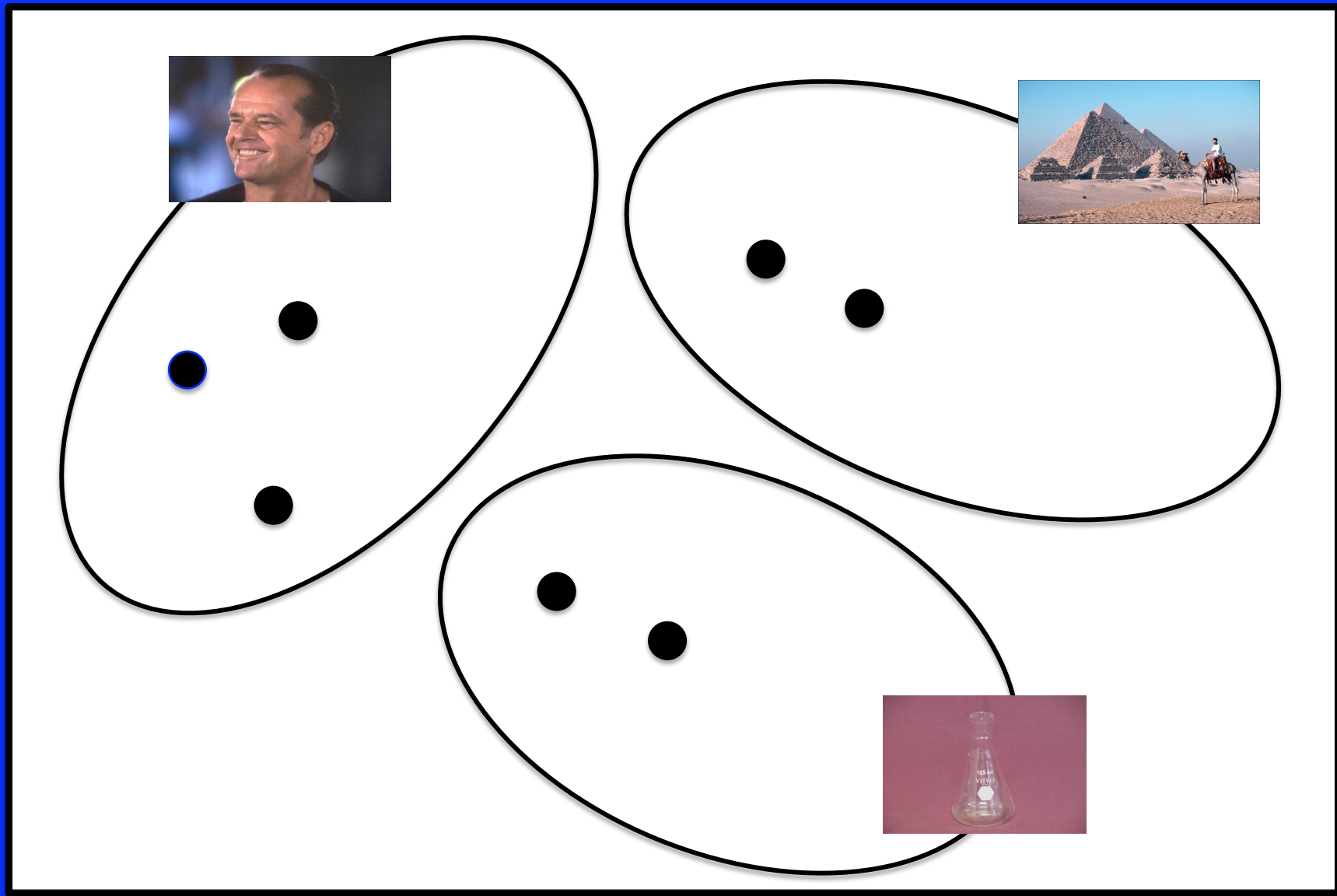
Context Space: Encoding



Context Space: Encoding



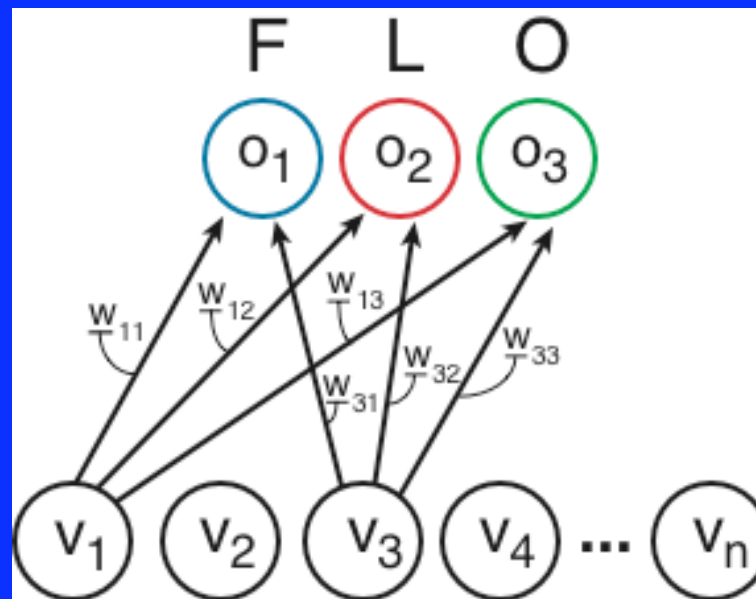
Context Space: Retrieval



Analysis strategy

- Part 1: Feed fMRI data from the study phase into a pattern classification algorithm
- Train the pattern classifier to recognize the brain patterns associated with studying faces vs. locations vs. objects

Neural network classifier



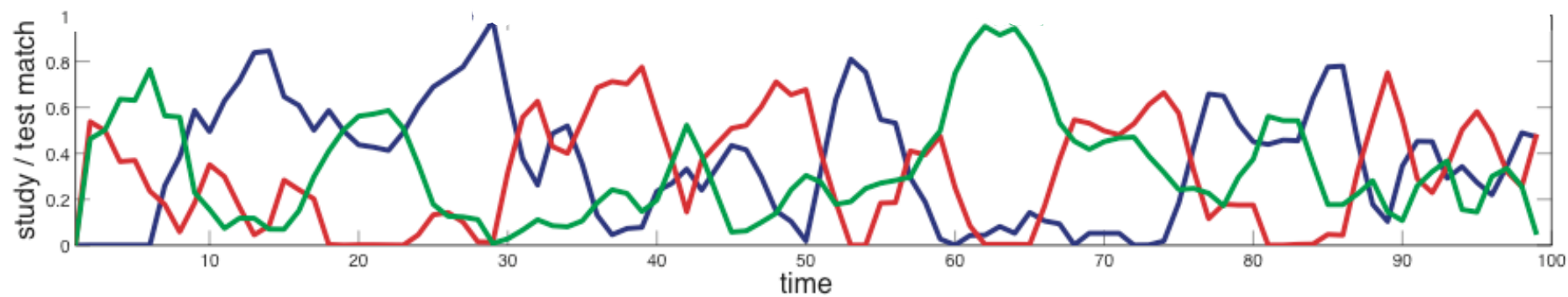
- Mapping from voxel activity values to output units (one per category)

Analysis strategy

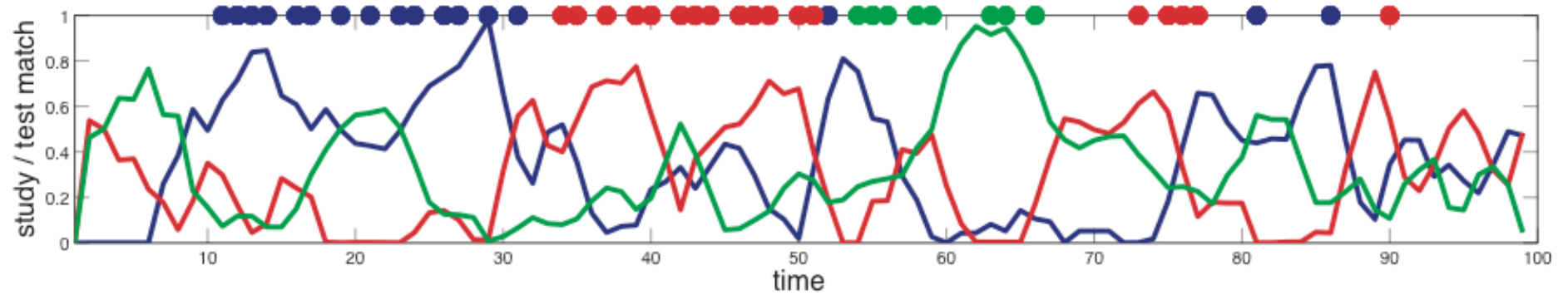
- Part 2: Apply the trained classifier to brain data from the retrieval phase
- Use the classifier to track, second-by-second, how well the subject's brain state at retrieval matches their brain state when they were studying faces vs. locations vs. objects

Predictions

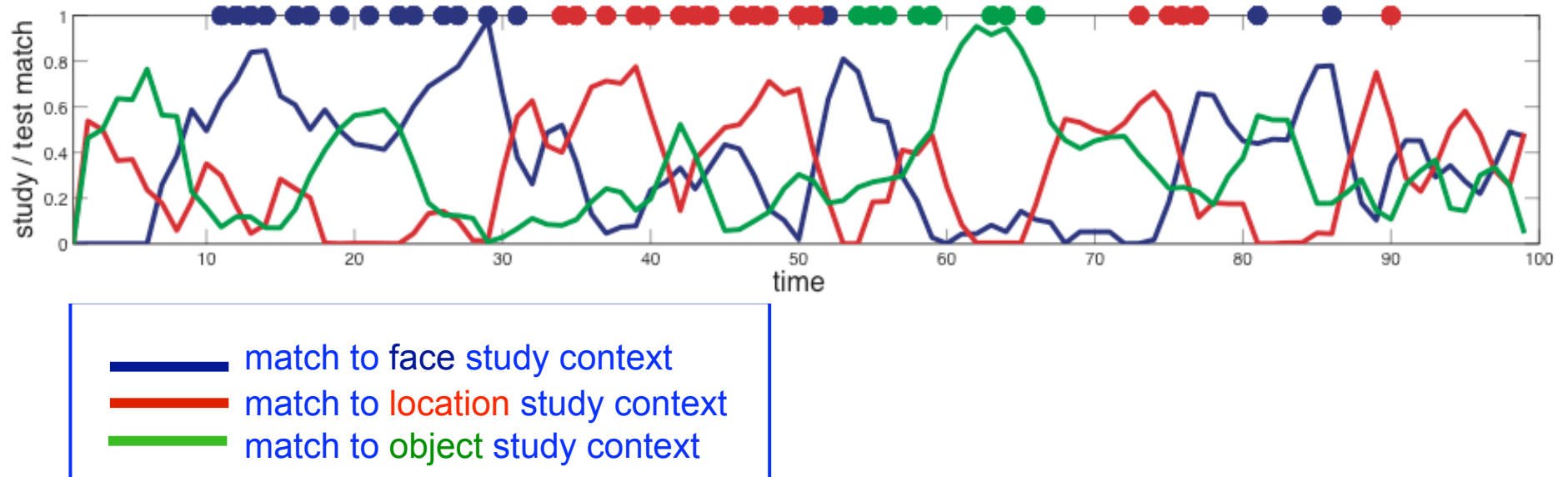
- As subjects try recall faces, locations, and objects, their brain state should come into alignment with the brain states associated with studying faces, locations, and objects
- This neural measure of category-specific contextual reinstatement should predict recall



- match to face study context
- match to location study context
- match to object study context



- match to face study context
- match to location study context
- match to object study context



- Reinstatement of category-specific brain activity correlated very strongly with recall behavior
- Category-specific brain activity started to emerge several seconds before subjects recalled items from that category
- We were able predict what category of item subjects would recall (with $>$ chance accuracy) based on data collected ~ 5 seconds **before** subjects recalled the item

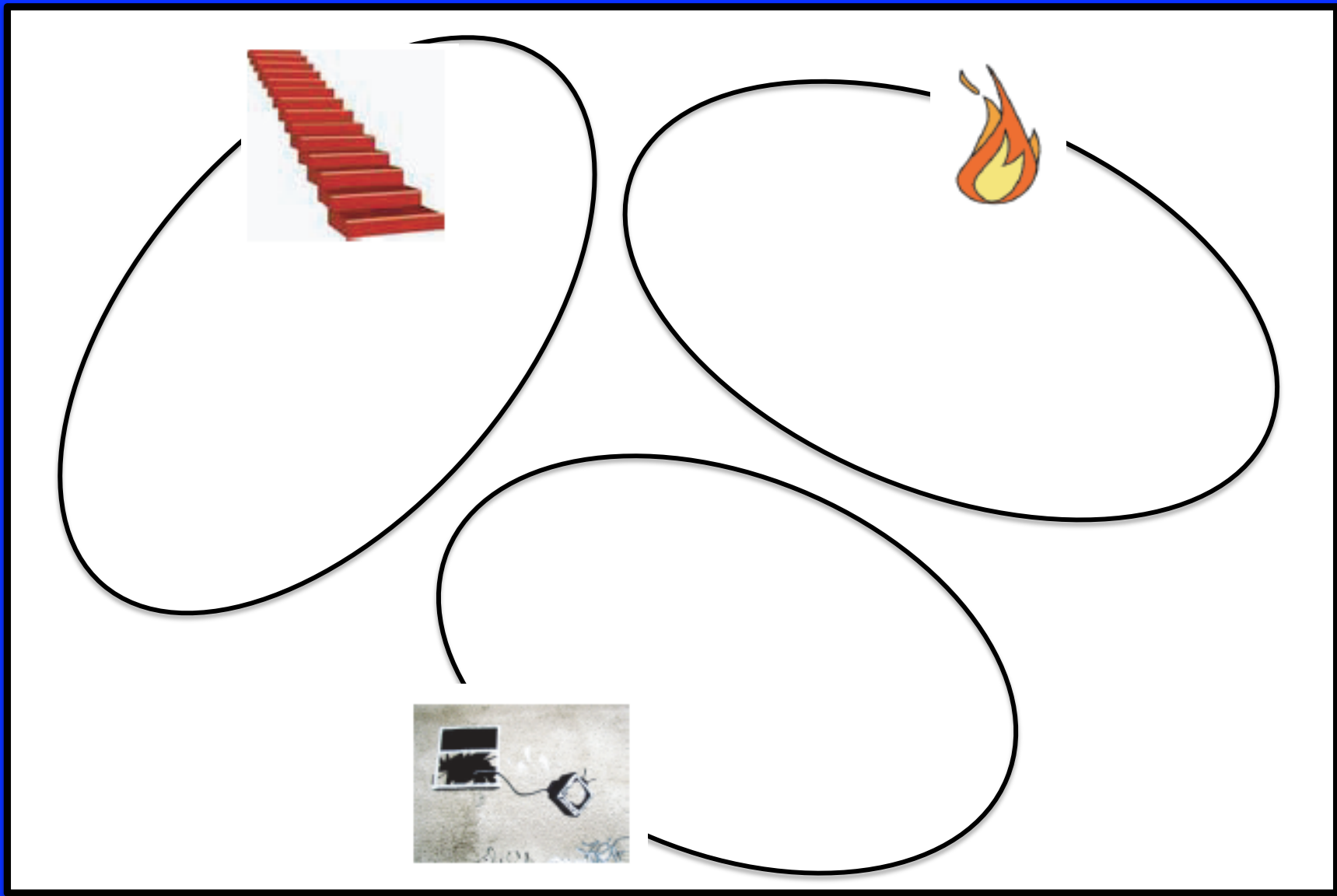
Shortcomings

- **Item** information was confounded with **context** information
- Face, location, & object activity at test may reflect subjects thinking about the **items** as opposed to subjects reinstating detailed “mental contexts” from the study phase
- Solution: Design a new experiment where items are arbitrarily assigned to contexts

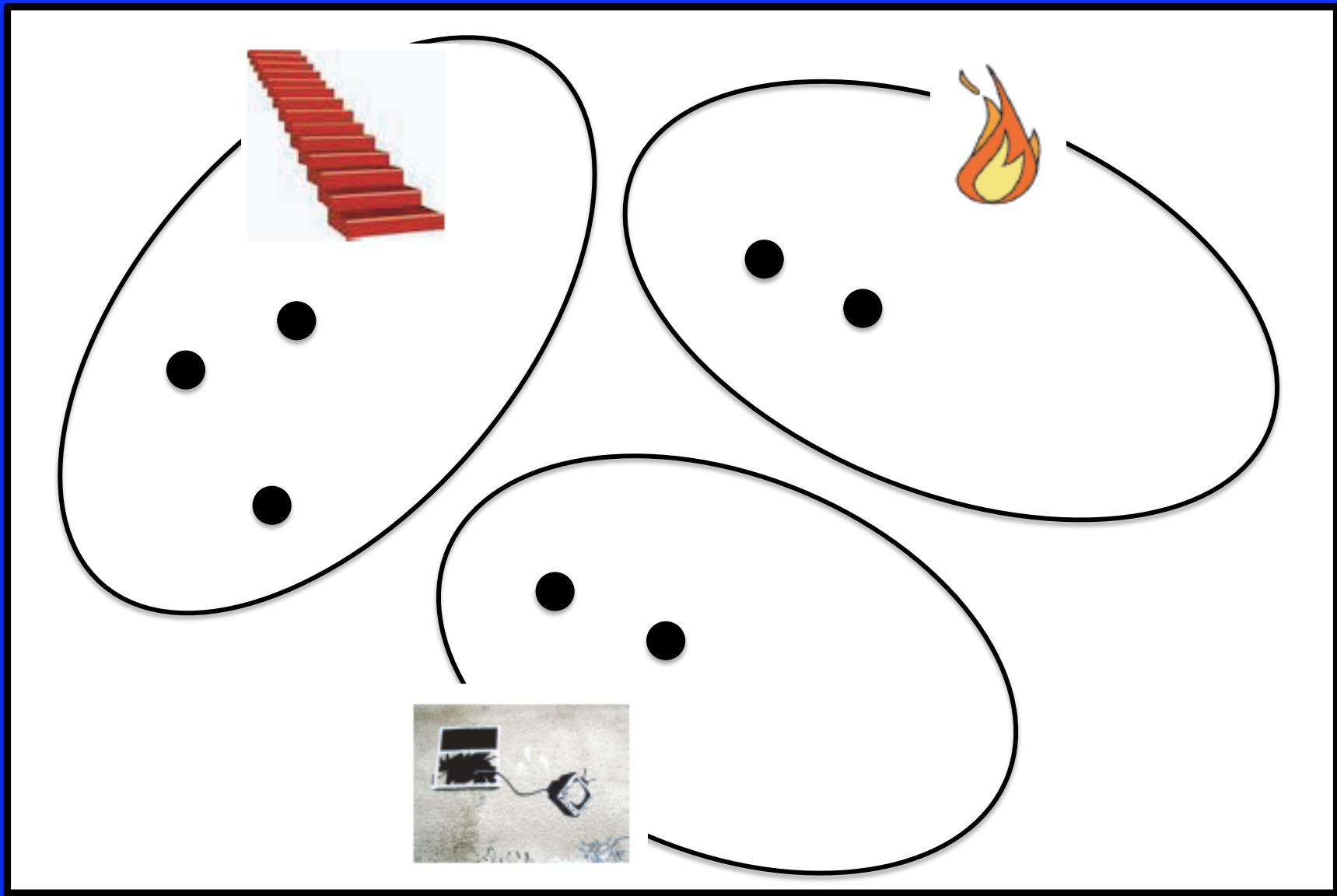
Bonfire study (Detre et al., 2007)

- Stimuli were concrete noun words
- Words were randomly assigned to one of 3 contexts:
 - Throw on bonfire
 - Carry up stairs
 - Drop out of window
- Train classifier (on study phase data) to recognize these three contexts
- Use the trained classifier to measure reinstatement of these contexts at test

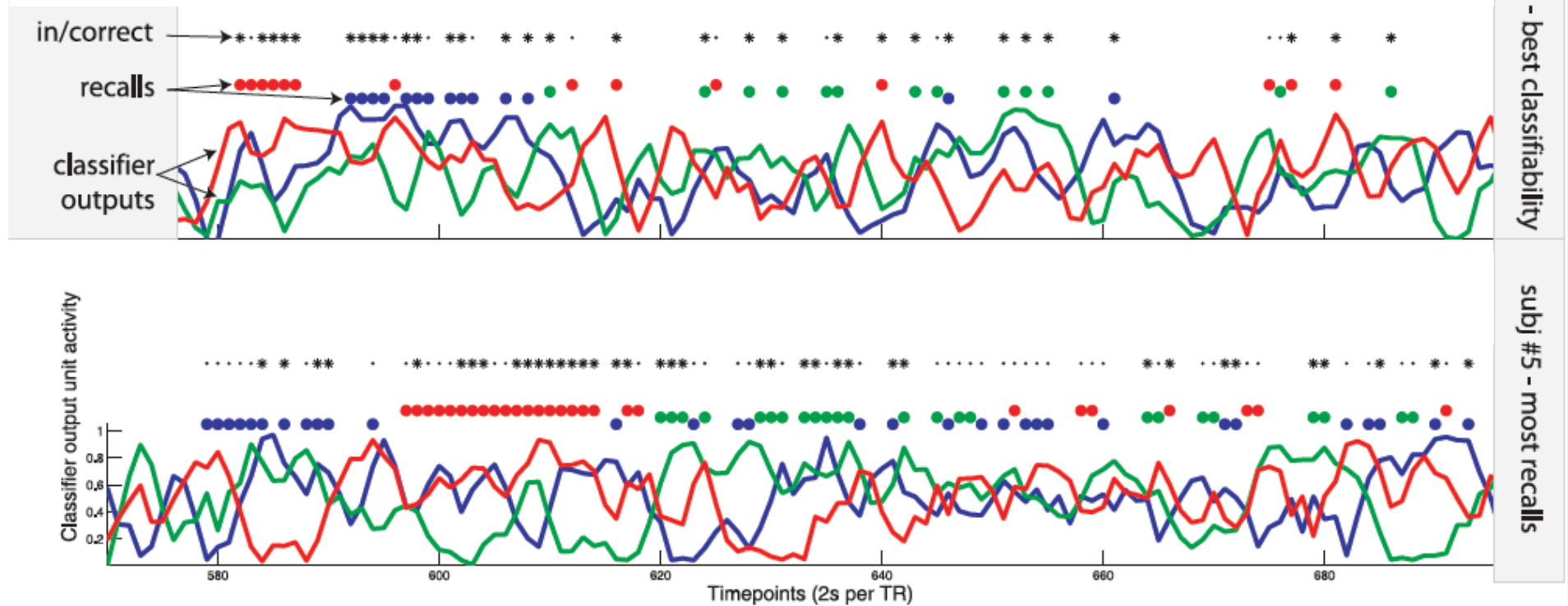
Context Space



Context Space



Study-recall generalization - 2 sample subjects - final free recall only



- Blue = bonfire, Red = stairs, Green = window
- Average percent correct across 8 subjects = 42%
(chance = 33%; range = 27% - 74%)

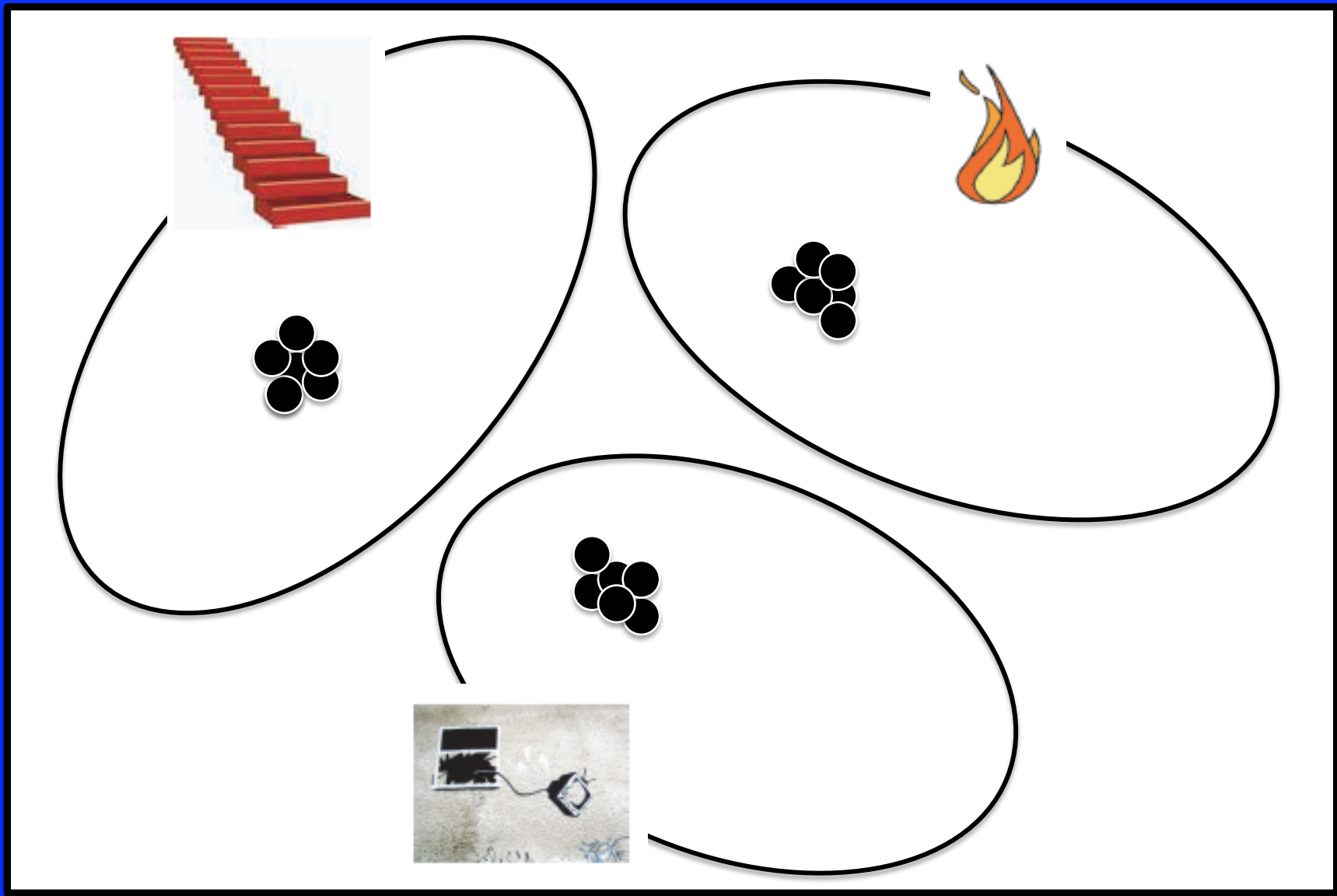
Bonfire study (Detre et al., 2007)

- Given that these results weren't so great, we decided to look more closely at study-phase data
- We ran a cross-validation analysis to assess whether the three contexts **elicit discriminable neural patterns at study**
- Study-phase cross-validation results were not too great either (47% accurate; chance = 33%)
- What might be responsible for these less-than-great results?

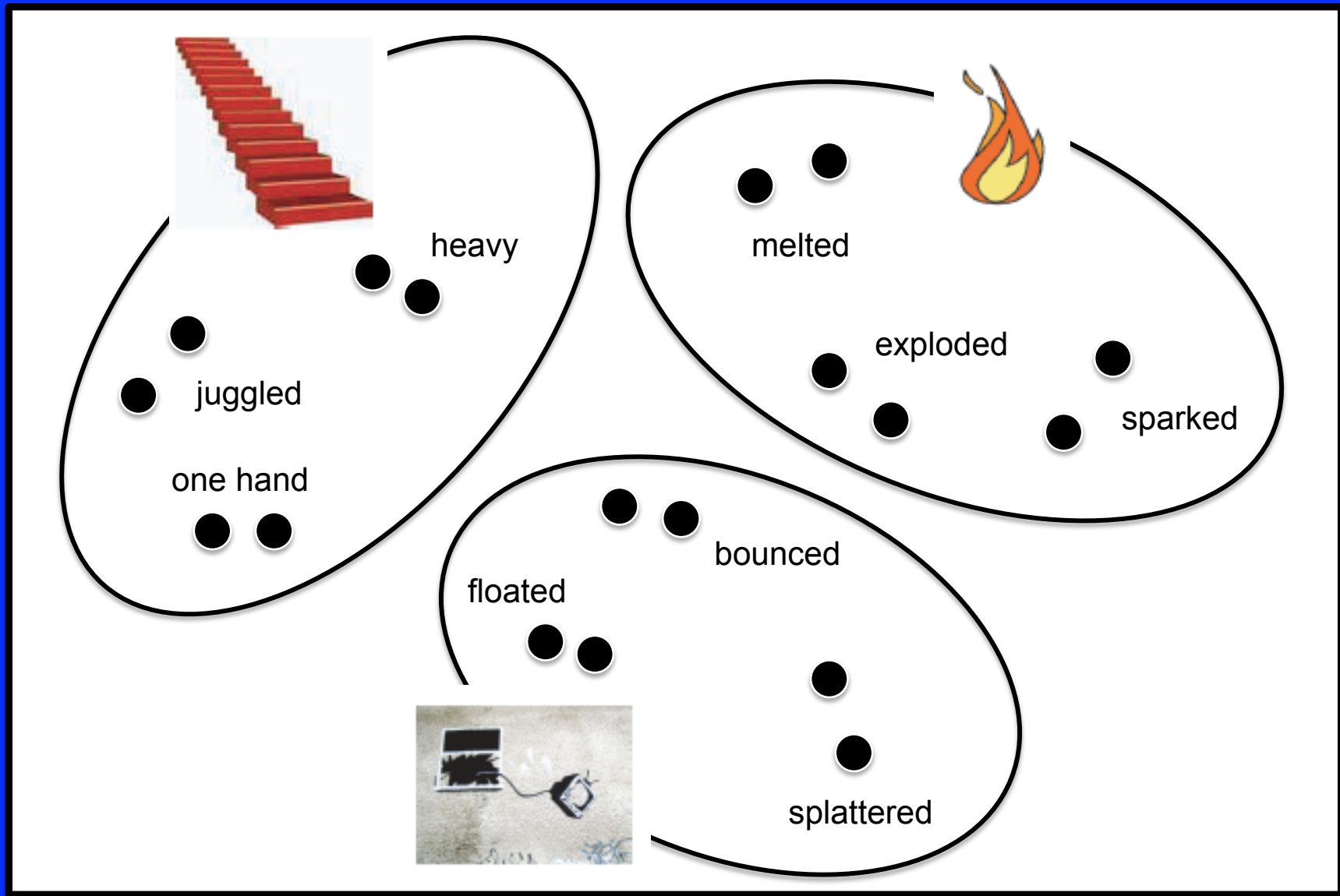
Bonfire study (Detre et al., 2007)

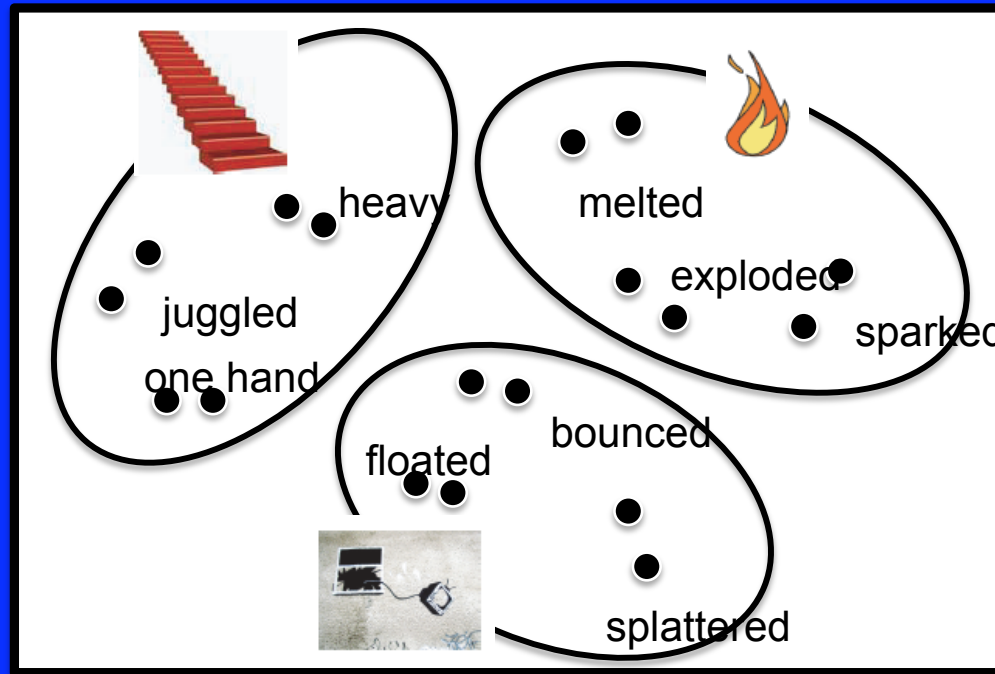
- Tradeoff between classifiability and memory performance
- To maximize classifier performance, representations should be **consistent**
- However, if you always think about the bonfire in exactly the same way, the bonfire context cue will become **overloaded**, leading to poor memory performance

Context Space

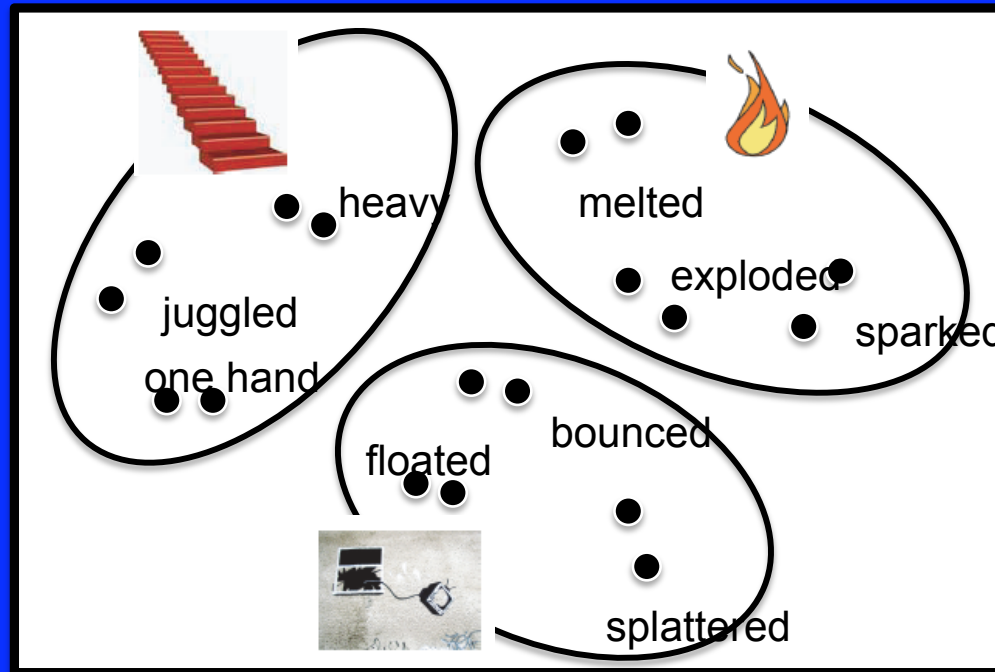


Context Space

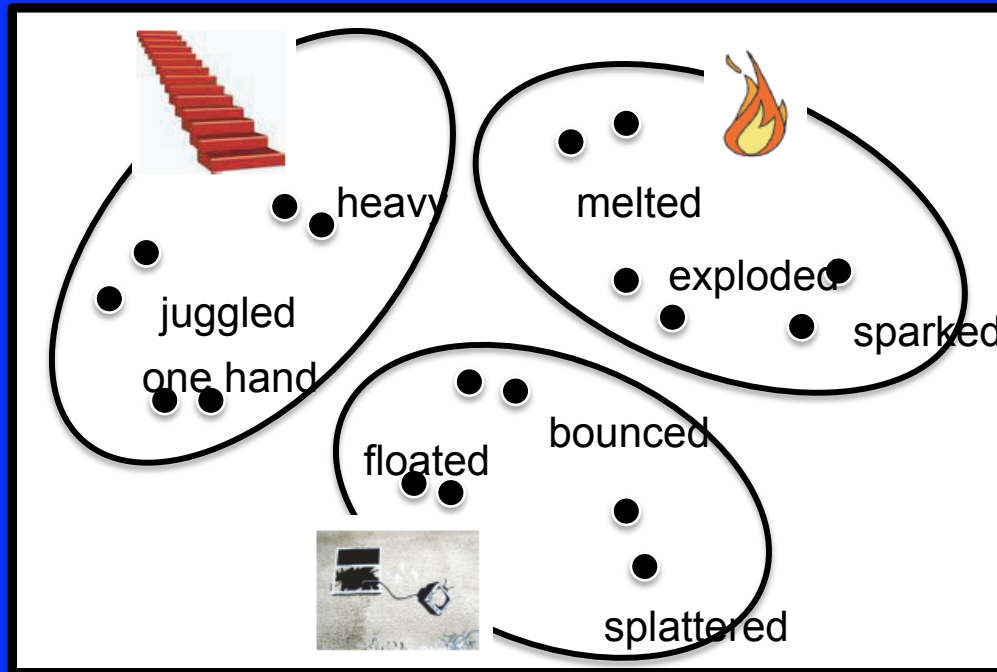




- This heterogeneity hurts classification but helps memory (by preventing cue overload)
- It might be possible to classify bonfire vs. stairs vs. window if we had more training data
- However, is this really worthwhile?



- We have established that “bonfire”, “stairs”, and “window” by themselves are overloaded memory cues
- To know what people are going to recall, we need to know more precisely where subjects are (mentally) in this space



- One possibility: Train the classifier to recognize more points in the space
- Ask subjects to perform specific sub-types of encoding within a context (e.g., STAIRS – juggle this; BONFIRE – melt this)
- Train the classifier on these encoding sub-types
- This still doesn't solve the problem!

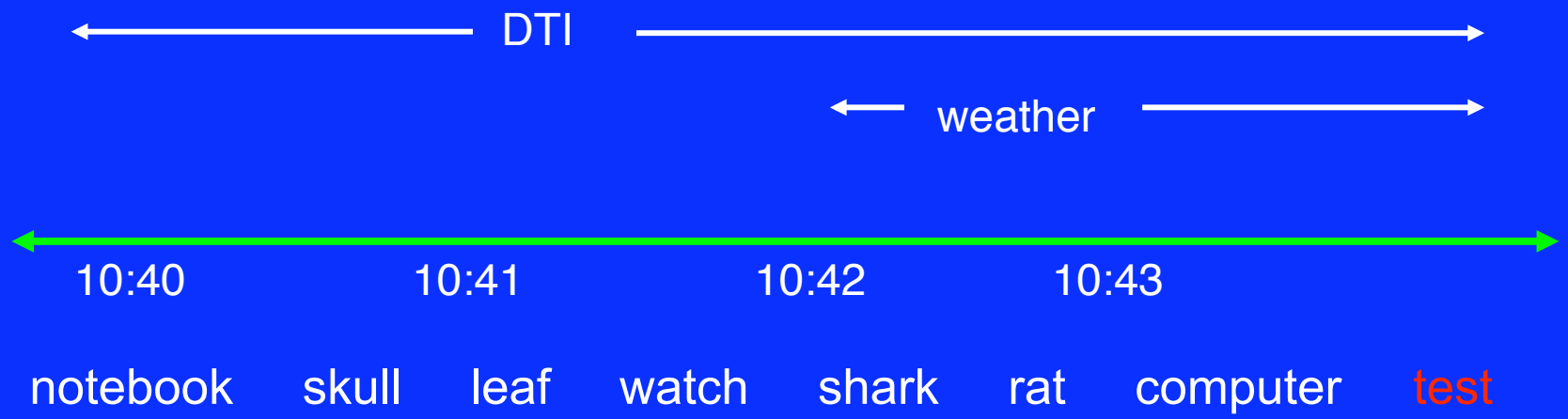
Beyond Classification

- Key claim of context models: Item representations are linked to **all other active thoughts**
- So what we really want is an efficient representation of the subject's **entire cognitive state**
- Naïve approach: Use the whole-brain activity vector as a proxy for the “context vector”
- Issue: Not all variance in the BOLD signal is cognitively relevant
- If we could isolate the “cognitively relevant” part of the whole-brain activity pattern, this might be a useful context representation

Beyond Classification

- We should also be able to use behavioral data **on context shift effects** to constrain the process of finding the neural context vector
- Numerous studies have explored how **interposing mental activities** during the study phase affects memory

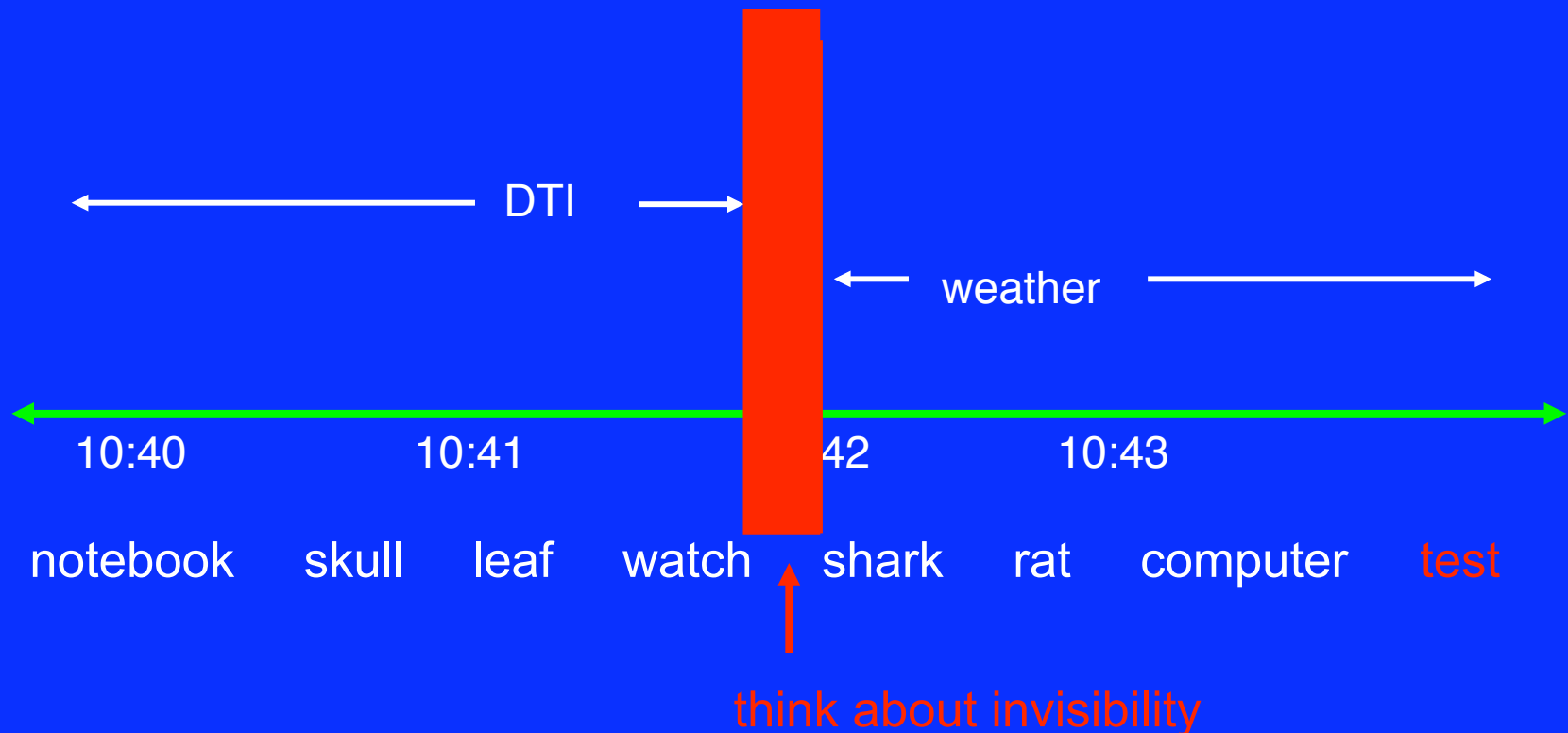
Contextual Disruption



Contextual Disruption



Contextual Disruption



- Activities that **strongly** disrupt context should impair recall of items studied **prior** to that activity and improve recall of items studied **after** that activity (Sahakyan & Kelley, 2002)

Beyond Classification

- We can use experiments like this to **rank mental activities** in terms of how much they disrupt recall, and we can use this to **infer** how much these activities disrupt context
- Key desiderata for neural context vector:
- Behavioral manipulations that are known to have a large effect on context (e.g., “think about invisibility”) should have a large effect on the neural context vector
- Behavioral manipulations that are known to have a small effect on context should have a small effect on the neural context vector

Focus on MTL

- The medial temporal lobes actually **do** the binding of item and context
- Thus, context information needs to be represented in MTL
- Instead of looking at the whole brain, it should be possible to do high-resolution imaging of MTL
- Use the MTL pattern as the “context vector”

Gradual Changes in Hippocampal Activity Support Remembering the Order of Events

Joseph R. Manns,^{1,3} Marc W. Howard,² and Howard Eichenbaum^{1,*}

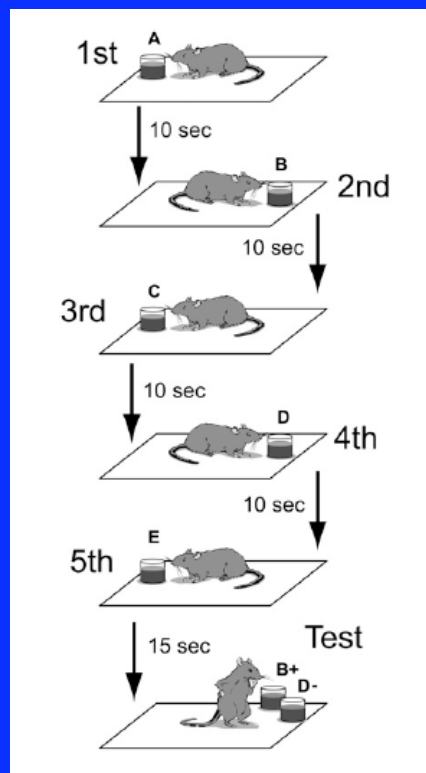
¹Center for Memory and Brain, Boston University, Boston, MA 02215, USA

²Department of Psychology, Syracuse University, Syracuse, NY 13244, USA

³Present address: Department of Psychology, Emory University, Atlanta, GA 30322, USA.

*Correspondence: hbe@bu.edu

DOI 10.1016/j.neuron.2007.08.017



- Show rats a series of odors
- Train rats to perform recency judgments
- Record multi-unit activity from CA1

Gradual Changes in Hippocampal Activity Support Remembering the Order of Events

Joseph R. Manns,^{1,3} Marc W. Howard,² and Howard Eichenbaum^{1,*}

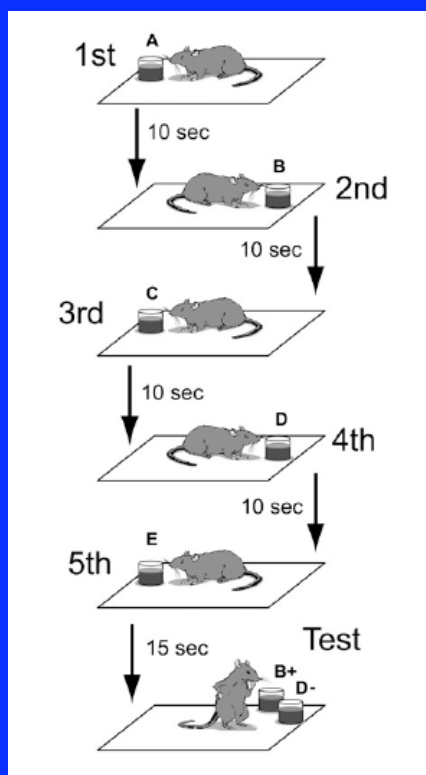
¹Center for Memory and Brain, Boston University, Boston, MA 02215, USA

²Department of Psychology, Syracuse University, Syracuse, NY 13244, USA

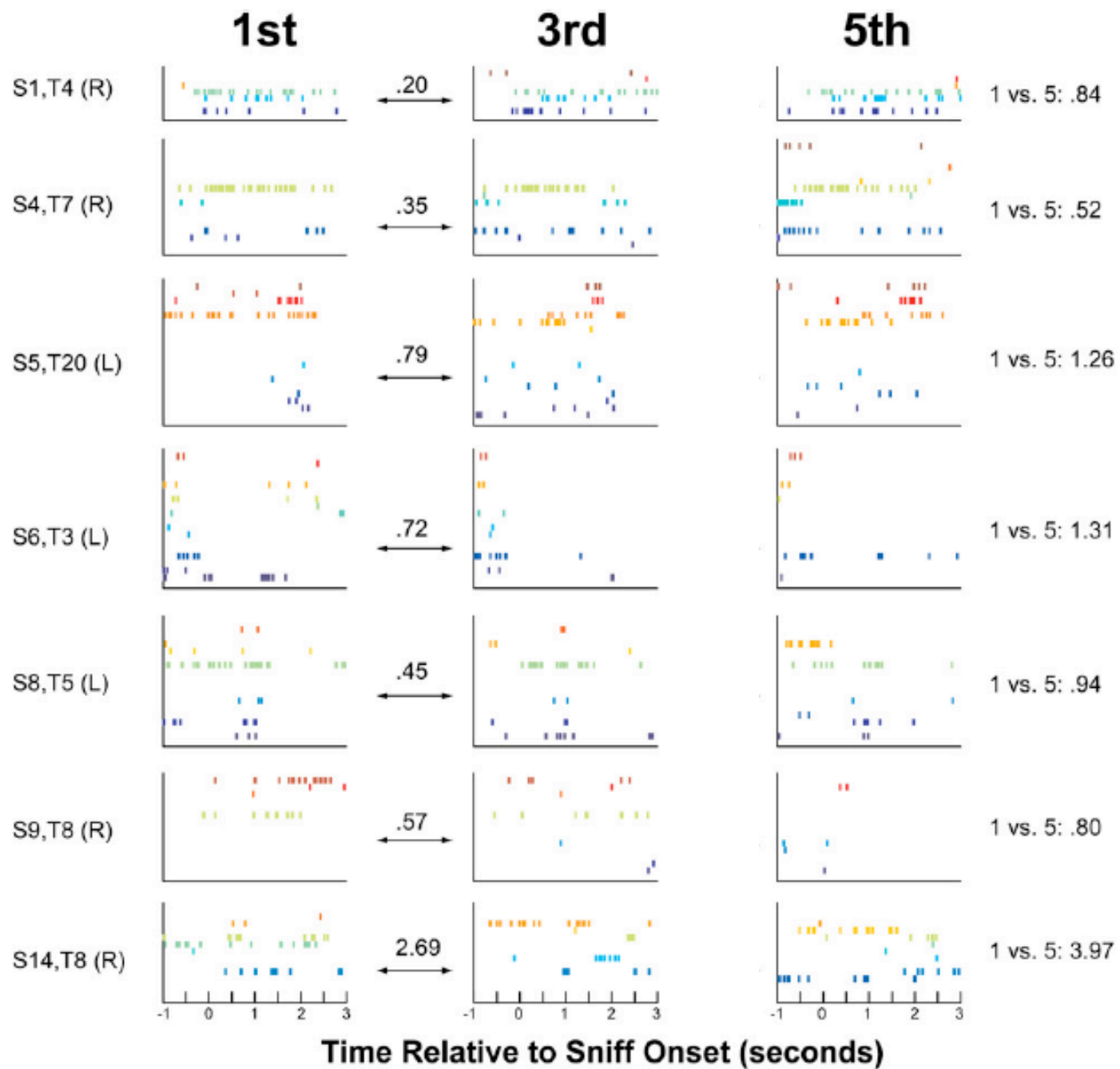
³Present address: Department of Psychology, Emory University, Atlanta, GA 30322, USA.

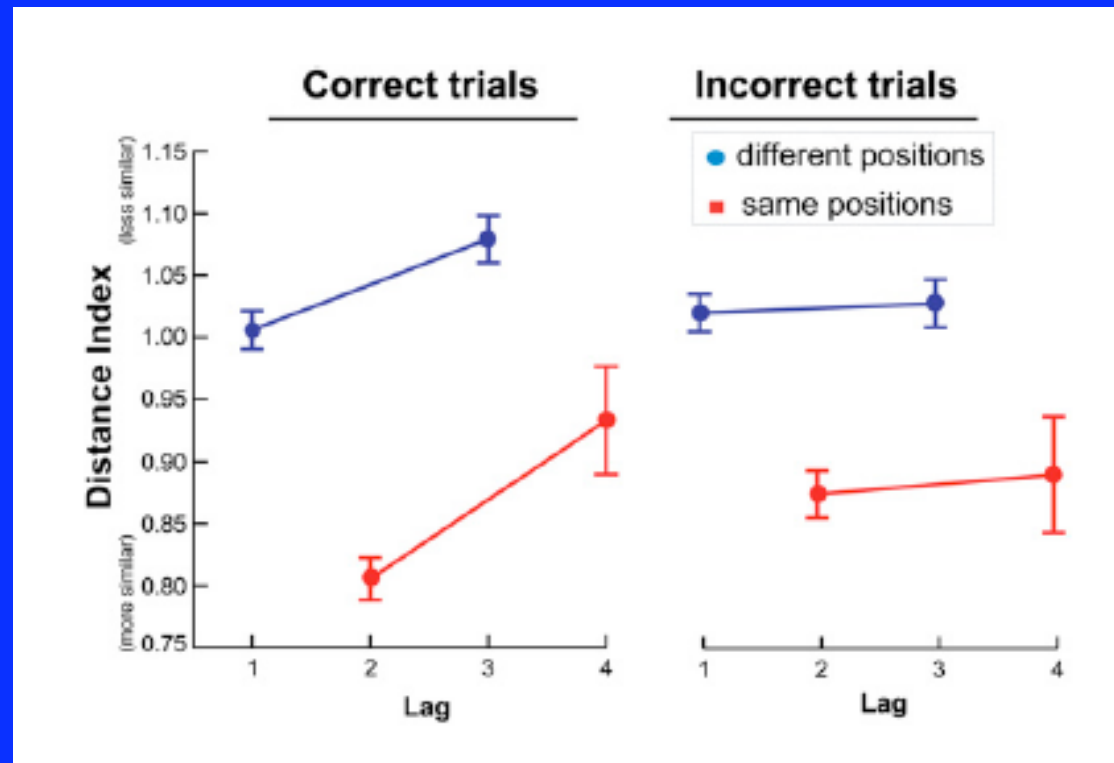
*Correspondence: hbe@bu.edu

DOI 10.1016/j.neuron.2007.08.017



- Use this multi-unit CA1 recording as a neural context vector
- Measure how much the context vector drifts during the encoding phase
- Use this to predict accuracy
- Intuitively: The more the context vector drifts between items, the more temporally discriminable the items will be



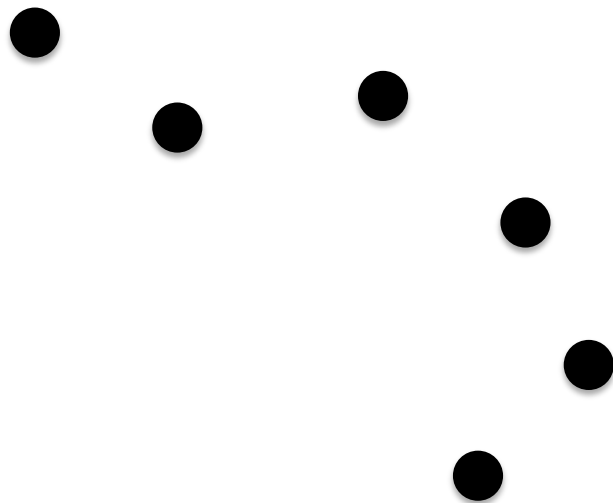


- It works: Increased “contextual drift” between items predicts increased accuracy

Memory Search: Summary

- Tremendous scientific payoff if we can image how subjects' mental context evolves during encoding and retrieval

Context Space



Memory Search: Summary

- Classification methods provide some insight into memory search...
- but the amount of information that we can glean is limited
- and the information that we get using classification methods is **not specific enough** to test our (very specific) mathematical models of memory search

Overall Summary

- Classifier methods can be used to track time-varying cognitive states
- Train on well-defined cognitive states, generalize to messy cognitive states
 - Negative priming: Train on target, generalize to distractor
 - Free recall: Train on study, generalize to test

Overall Summary

- The problem that got my lab into the classification business – memory search – has proved to be interestingly resistant to standard classification methods
 - We need to track subjects' position in a very high dimensional mental space
 - High-dimensional is not a problem, if the dimensions are well-defined
 - Mitchell et al. (2008) were able to decode what word subjects were thinking of, by representing each word in a 25-dimensional “semantic feature space”, and then learning the brain patterns associated with each dimension

Overall Summary

- Kay et al. (2008) were able to decode what photo subjects were thinking of, by representing photos in terms of low-level visual features, and then learning the brain patterns associated with these low-level features
- This approach is harder to apply to memory search, because the dimensions of the contextual “search space” are not always apparent beforehand



Princeton Computational Memory Lab



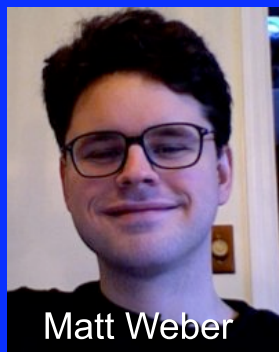
Ehren Newman



Greg Detre



Ken Norman



Matt Weber



Susan Robison



Chris Moore



Adler Perotte



Joel Quamme



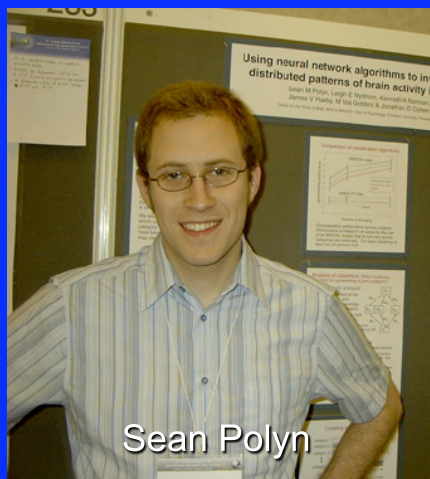
Per Sederberg



David Weiss



Princeton Computational Memory Lab





The fMRI analyses were run using the
Princeton Multi-Voxel Pattern Analysis Toolkit
downloadable from:
<http://www.csbmb.princeton.edu/mvpa>