

# LING5702: Lecture Notes 8

## A Model of Ambiguity in Sentence Processing

We have seen how complex ideas can be encoded and decoded into sentences.

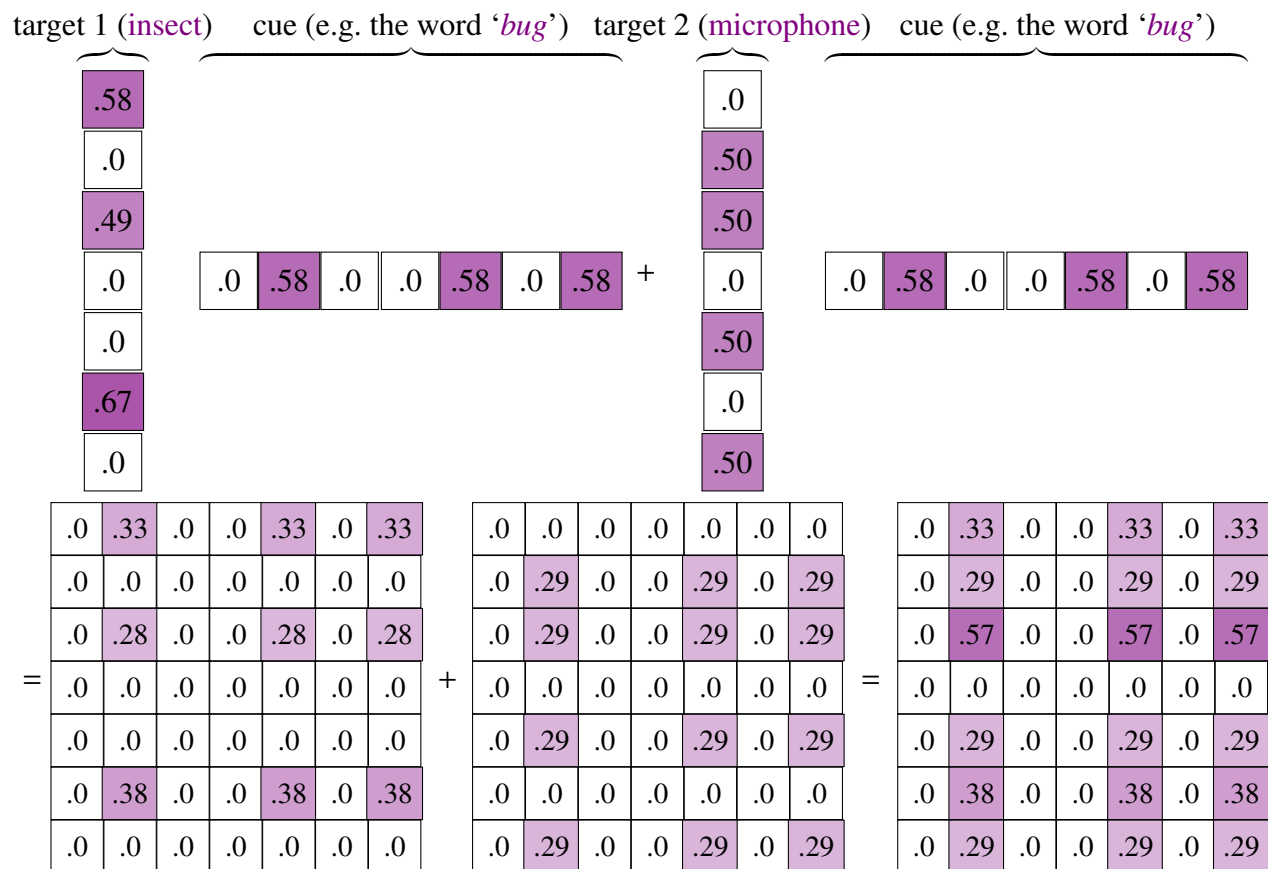
This lecture will describe how this process controls for ambiguity.

### Contents

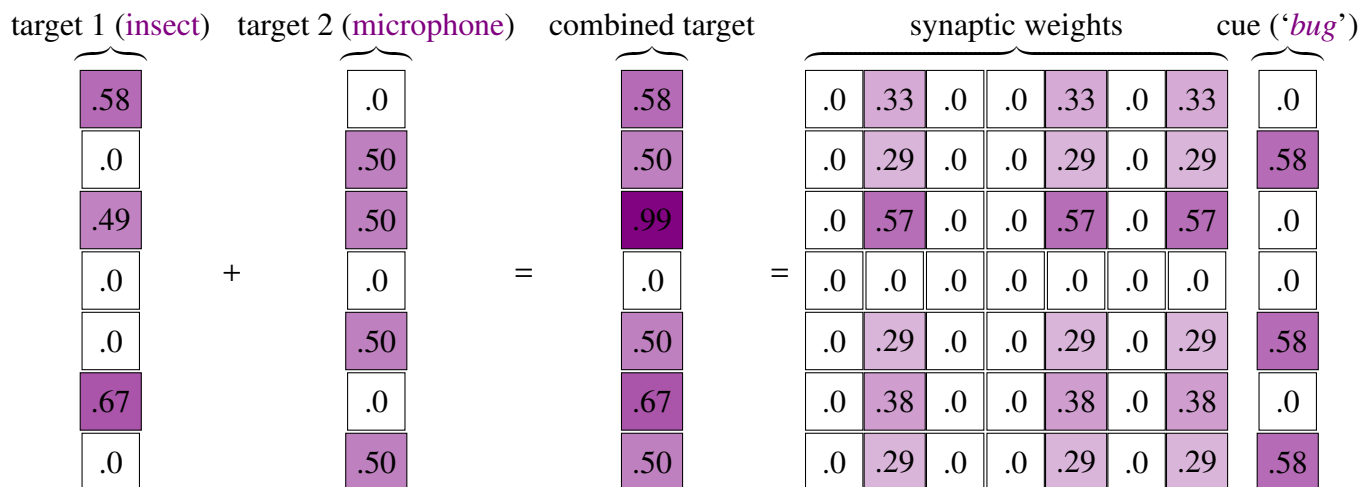
8.1	Ambiguity as superposition [Smolensky, 1990, Rasmussen & Schuler, 2018] . . . .	1
8.2	Propagation of ambiguity . . . . .	3
8.3	Resolution of ambiguity . . . . .	5
8.4	Probabilistic disambiguation [Friston, 2010] . . . . .	7
8.5	Neural ambiguity resolution . . . . .	9

### 8.1 Ambiguity as superposition [Smolensky, 1990, Rasmussen & Schuler, 2018]

If multiple targets are associated with the same cue vector:

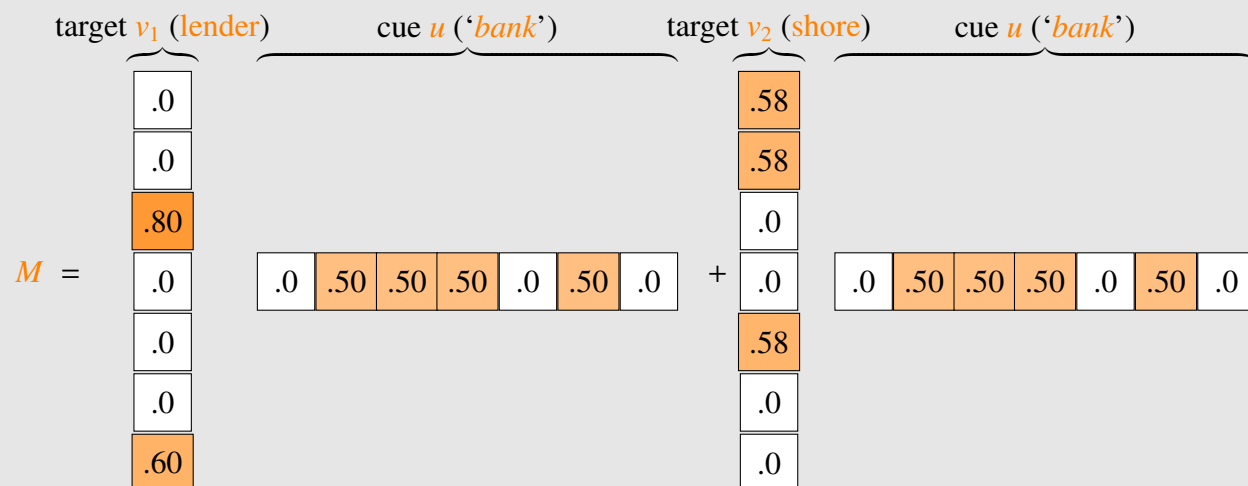


when associative memory is cued using that vector, the result is both vectors, **superposed**:

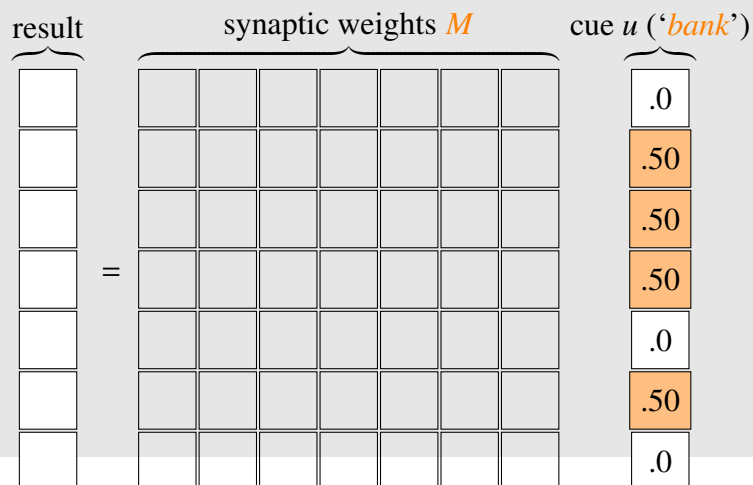


### Practice 8.1:

If associative memory  $M$  is made from one cue  $u$  and two targets  $v_1$  and  $v_2$ , as below:



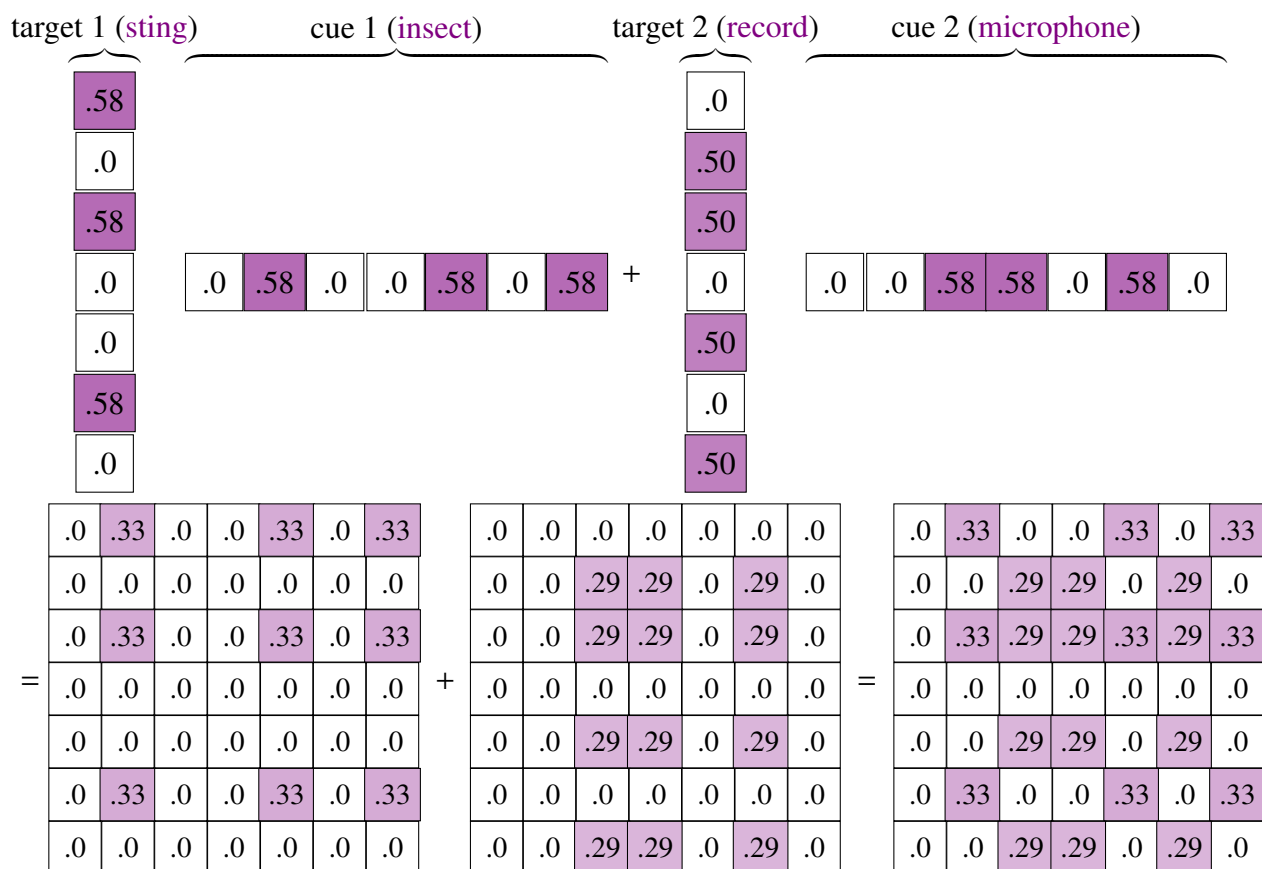
what is the result of cueing  $M$  with  $u$ ? (HINT: you don't have to calculate the matrix!)



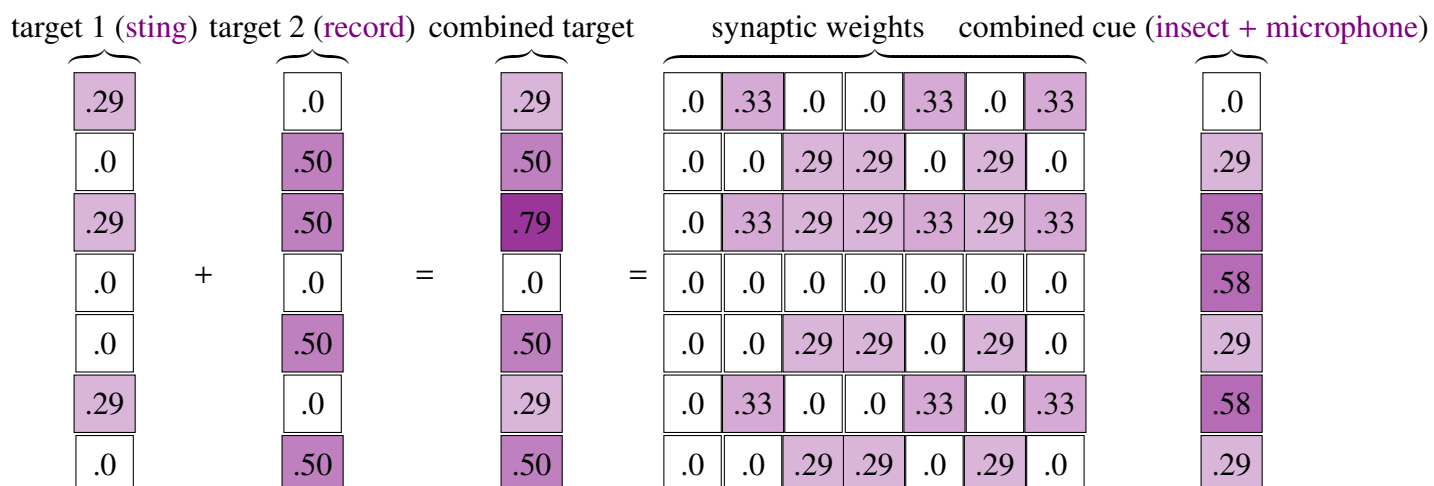
Describe the result in terms of  $v_1$  and  $v_2$ .

## 8.2 Propagation of ambiguity

If we take a (non-interfering) set of associations:



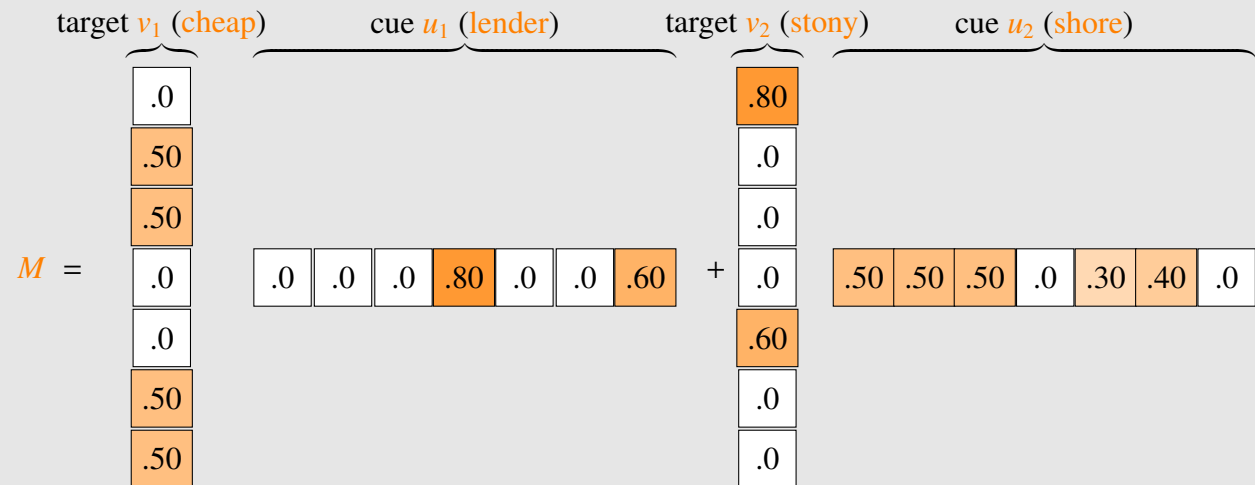
and cue them with a combination of states, we get a proportional combination of targets:



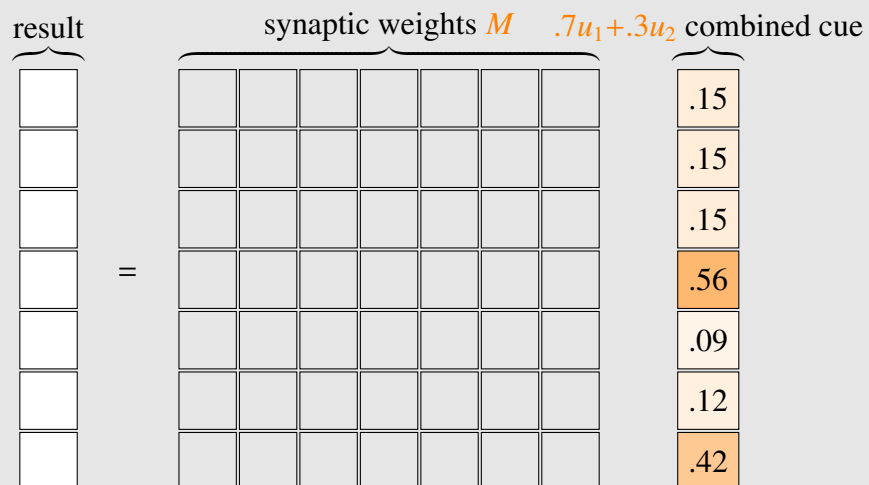
This is important because it allows ambiguity to propagate through a mental process.

### Practice 8.2:

If associative memory  $M$  is made from cues  $u_1$  and  $u_2$  and targets  $v_1$  and  $v_2$ , as below:



what is the result of cueing  $M$  with a mixture of  $.7u_1$  and  $.3u_2$ ? (HINT: don't calculate the matrix!)



Describe the result in terms of  $v_1$  and  $v_2$ .

### 8.3 Resolution of ambiguity

Recall the outer product of two vectors produces a matrix with pointwise products:

$$\begin{array}{c} \text{target } v_1 \end{array} \begin{array}{c} \text{cue } u_1 \end{array} = \begin{array}{c} \text{associations } M_1 \end{array} \quad (1)$$

Diagram illustrating the outer product of two vectors to produce a matrix:

- target  $v_1$** : A vertical vector with values  $.0, .71, .0, .0, .71$ .
- cue  $u_1$** : A horizontal vector with values  $.0, .71, .0, .0, .71$ .
- associations  $M_1$** : A 5x5 matrix resulting from the pointwise products of the target and cue vectors. The matrix is:
 

.0	.0	.0	.0	.0
.0	.50	.0	.0	.50
.0	.0	.0	.0	.0
.0	.0	.0	.0	.0
.0	.50	.0	.0	.50

This generalizes to triples of vectors as a **tensor product**:

$$\begin{array}{c} \text{target } v'_1 \end{array} \begin{array}{c} \text{cue } v_1 \end{array} \begin{array}{c} \text{cue } u_1 \end{array} = \begin{array}{c} \text{tensor } \mathcal{T}_1 \end{array}$$

Diagram illustrating the tensor product of three vectors to produce a 3D tensor:

- target  $v'_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .
- cue  $v_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .
- cue  $u_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .
- tensor  $\mathcal{T}_1$** : A 5x5x5 3D tensor resulting from the pointwise products of the three vectors. The tensor is represented as a stack of 5x5 matrices, with the top matrix being:
 

.19	.0	.0	.19	.19
.0	.0	.0	.0	.0
.0	.0	.0	.0	.0
.19	.0	.0	.19	.19
.19	.0	.0	.19	.19

Targets are then cued by multiplication (left-associative):

$$\begin{array}{c} \text{target } v'_1 \end{array} \left( \begin{array}{c} \text{tensor } \mathcal{T}_1 \end{array} \begin{array}{c} \text{cue } u_1 \end{array} \right) \begin{array}{c} \text{cue } v_1 \end{array}$$

Diagram illustrating the left-associative multiplication of the target vector with the tensor and cue vector:

- target  $v'_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .
- tensor  $\mathcal{T}_1$** : A 5x5x5 3D tensor (as defined above).
- cue  $u_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .
- cue  $v_1$** : A vertical vector with values  $.58, .0, .0, .58, .58$ .

$$\begin{array}{c}
\text{target } v'_1 \\
\begin{array}{|c|} \hline .58 \\ \hline .0 \\ \hline .0 \\ \hline .58 \\ \hline .58 \\ \hline \end{array}
\end{array}
=
\begin{array}{c}
\text{matrix } \mathcal{T}_1 u_1 \\
\begin{array}{|c|c|c|c|c|c|} \hline .33 & .0 & .0 & .33 & .33 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .33 & .0 & .0 & .33 & .33 \\ \hline .33 & .0 & .0 & .33 & .33 \\ \hline \end{array}
\end{array}
\begin{array}{c}
\text{cue } v_1 \\
\begin{array}{|c|} \hline .58 \\ \hline .0 \\ \hline .0 \\ \hline .58 \\ \hline .58 \\ \hline \end{array}
\end{array}$$

This can be implemented with ‘switched’ connections.

This gives us a means to combine ‘top-down’ predictions with ‘bottom-up’ observations...

Build auto-associations of all states:

$$\begin{array}{c}
\text{result } v'_1 \text{ predicted } v_1 \text{ observed } u_1 \text{ result } v'_3 \text{ predicted } v_3 \text{ observed } u_3 \\
\text{(sting)} \quad \text{(sting)} \quad \text{('sting')} \quad \text{(sink)} \quad \text{(sink)} \quad \text{('sink')} \\
\begin{array}{|c|} \hline .71 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline .71 \\ \hline \end{array}
\otimes
\begin{array}{|c|} \hline .71 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline .71 \\ \hline \end{array}
\otimes
\begin{array}{|c|} \hline .71 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline .71 \\ \hline \end{array}
+
\begin{array}{|c|} \hline .0 \\ \hline .99 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline \end{array}
\otimes
\begin{array}{|c|} \hline .0 \\ \hline .99 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline \end{array}
\otimes
\begin{array}{|c|} \hline .0 \\ \hline .99 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline \end{array}
=
\begin{array}{c}
\text{tensor } \mathcal{T} \\
\begin{array}{|c|c|c|c|c|} \hline .35 & .0 & .0 & .0 & .35 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .35 & .0 & .0 & .0 & .35 \\ \hline \end{array}
\end{array}$$

Then cue on observed state to pick out compatible component of mixed source state:

$$\begin{array}{c}
\text{result } \approx v'_1 \\
\begin{array}{|c|} \hline .24 \\ \hline .0 \\ \hline .0 \\ \hline .0 \\ \hline .24 \\ \hline \end{array}
\end{array}
=
\begin{array}{c}
\text{tensor } \mathcal{T} \\
\begin{array}{|c|c|c|c|c|} \hline .35 & .0 & .0 & .0 & .35 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .0 & .0 & .0 & .0 & .0 \\ \hline .35 & .0 & .0 & .0 & .35 \\ \hline \end{array}
\end{array}
\begin{array}{c}
\text{observed } u_1 \\
\text{('sting')} \\
\begin{array}{|c|} \hline .71 \\ \hline .10 \\ \hline .0 \\ \hline .0 \\ \hline .71 \\ \hline \end{array}
\end{array}
\begin{array}{c}
\text{predicted state } v_1 + v_2 \\
\begin{array}{|c|} \hline .24 \\ \hline .0 \\ \hline .47 \\ \hline .47 \\ \hline .24 \\ \hline \end{array}
\end{array}$$

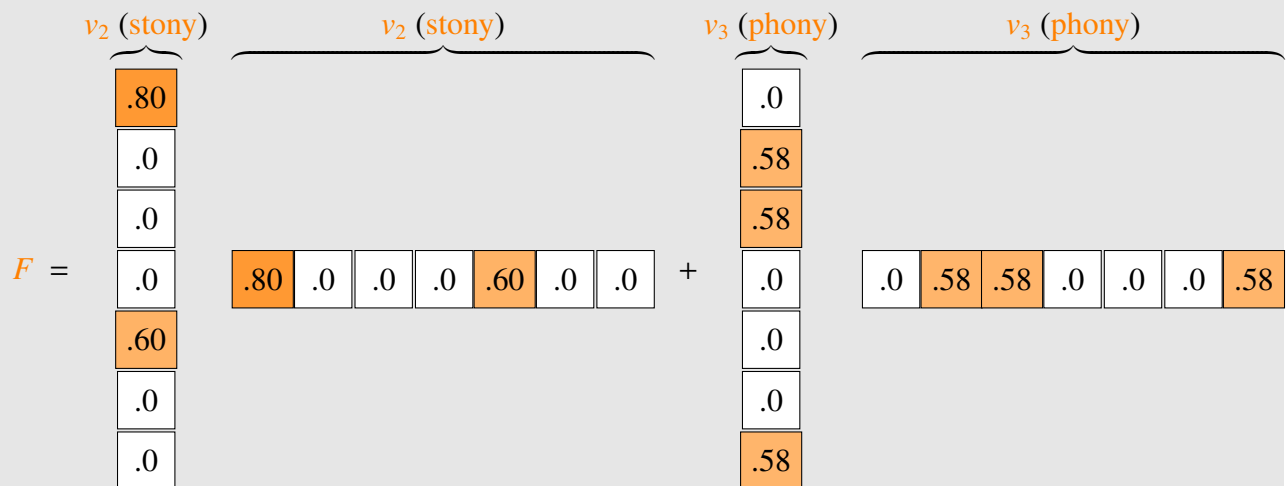
Note the magnitude of the target is reduced compared to the source.

This reduction correlates with reading time delays ('surprisal') on encountering unpredicted words.

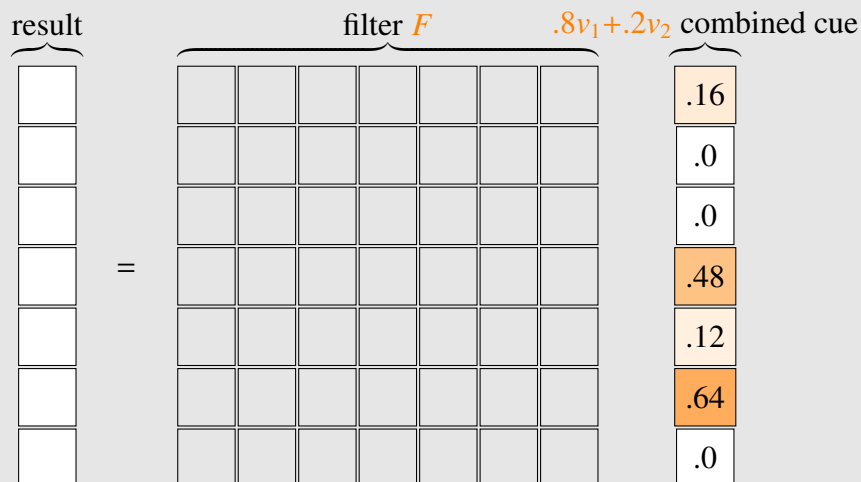
(It may take time proportional to the reduction to 'amp up' this state to unit magnitude.)

### Practice 8.3:

If a filter  $F$  is made from auto-associated vectors  $v_2$  and  $v_3$ , as below:



what is the result of cueing  $F$  with a mixture of  $.8v_1$  and  $.2v_2$ ? (HINT: don't calculate the matrix!)



## 8.4 Probabilistic disambiguation [Friston, 2010]

This multiplicative disambiguation is important because it allows probabilistic reasoning.

If:

- the prediction vector represents a **prior probability** of an idea  $P(i)$ , and
- the filter matrix represents a **likelihood**  $P(w | i)$  of a word  $w$  given an idea  $i$ ,

then:

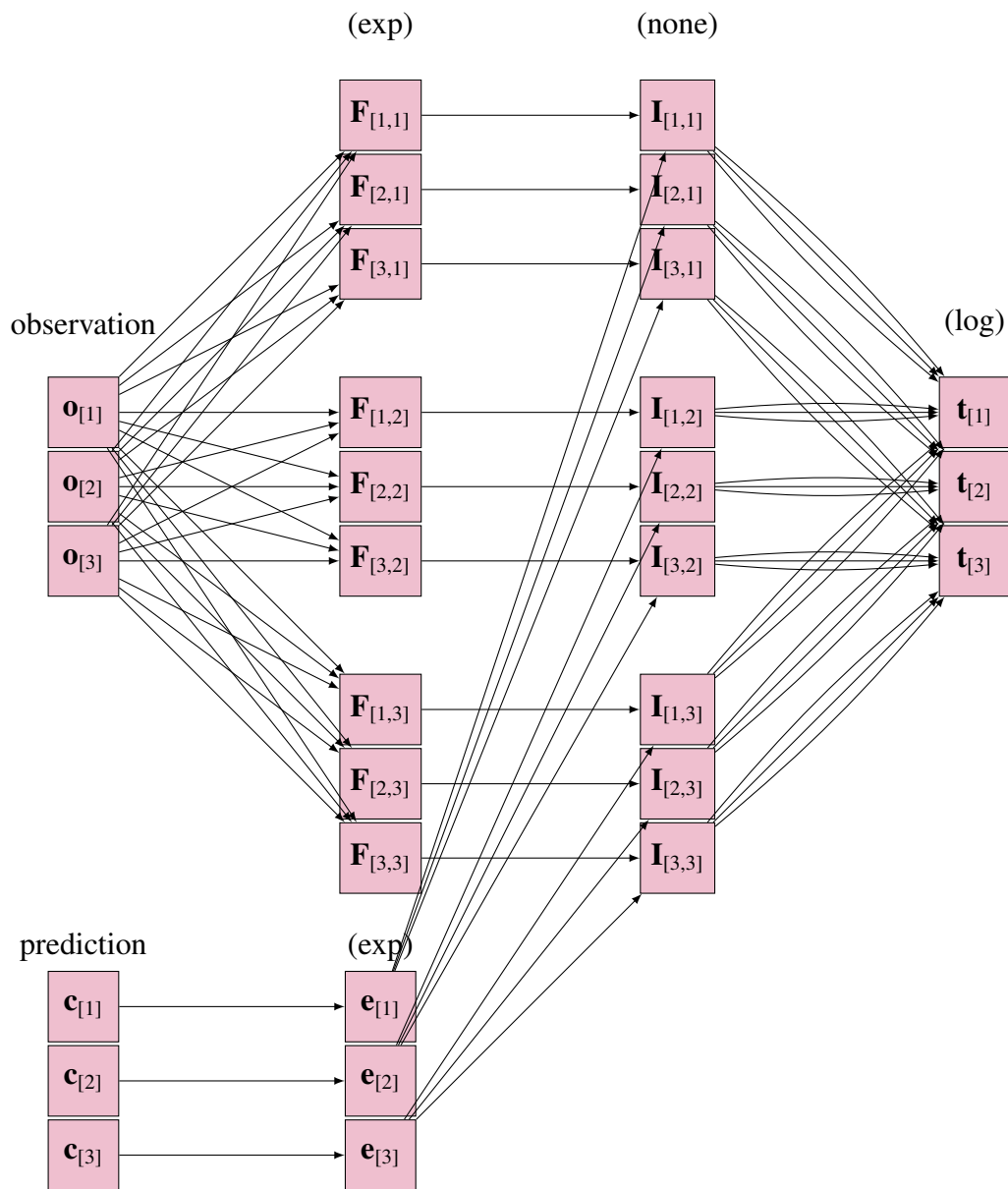
- the tensor model will calculate a correct **posterior probability** of the idea given the word!

$$P(i | w) = \frac{P(i) \cdot P(w | i)}{P(w)} \quad (\text{rescaled by a constant } P(w) \text{ to sum to one})$$



## 8.5 Neural ambiguity resolution

Here's a (hypothetical) neural implementation of tensor disambiguation:



1. observation  $\mathbf{o}$  is weighted by tensor weights to give filter  $\mathbf{F}$  with exponential transfer function.
2. filter  $\mathbf{F}$  is multiplied by prediction  $\mathbf{c}$  (adding exponentials) to give intermediate values  $\mathbf{I}$ .
3. intermediate values  $\mathbf{I}$  at several thresholds are added (log function) to get filtered prediction  $\mathbf{t}$ .

But any neural network can behave probabilistically if trained ‘autoregressively’ (to predict words).

## References

- [Friston, 2010] Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2), 127–138.
- [Rasmussen & Schuler, 2018] Rasmussen, N. E. & Schuler, W. (2018). Left-corner parsing with distributed associative memory produces surprisal and locality effects. *Cognitive Science*, 42(S4), 1009–1042.
- [Smolensky, 1990] Smolensky, P. (1990). Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artificial intelligence*, 46(1-2), 159–216.