

# **CN510: Principles and Methods of Cognitive and Neural Modeling**

## **Classical Conditioning and Outstar Network Outstar Learning Theorem**

### **Lecture 15**

Instructor: Anatoli Gorchetchnikov <[anatoli@bu.edu](mailto:anatoli@bu.edu)>

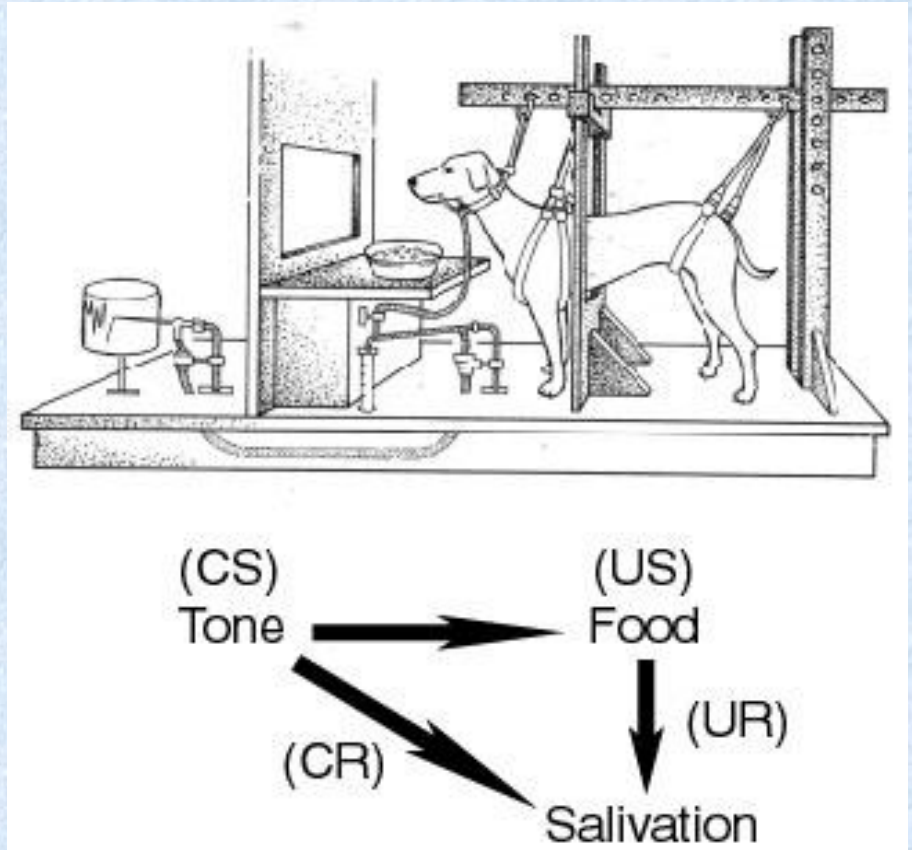
# Pavlovian Conditioning

Pavlov observed the result of repeated learning trials

Animal must make the following prediction:

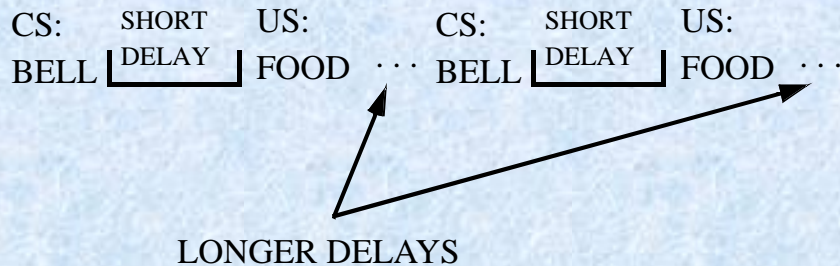
- When presented with a conditioned stimulus (CS), what, if any, unconditioned stimulus (US) will follow?

Prediction of US is evidenced by conditioned response (CR) that mimics the unconditioned response (UR)

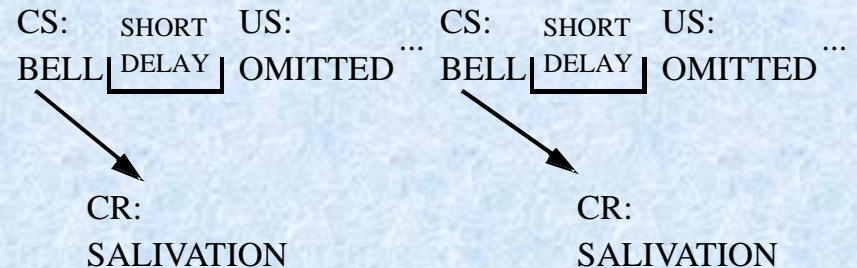


# Typical Classical Conditioning Paradigm

## (1) Learning / practice trials



## (3) Recall/performance trials

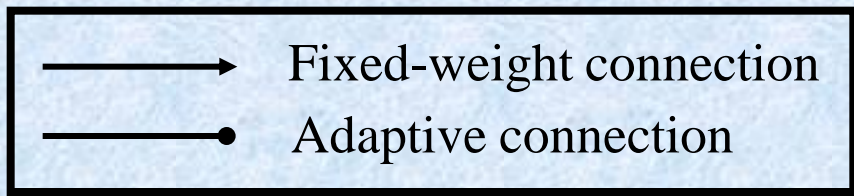
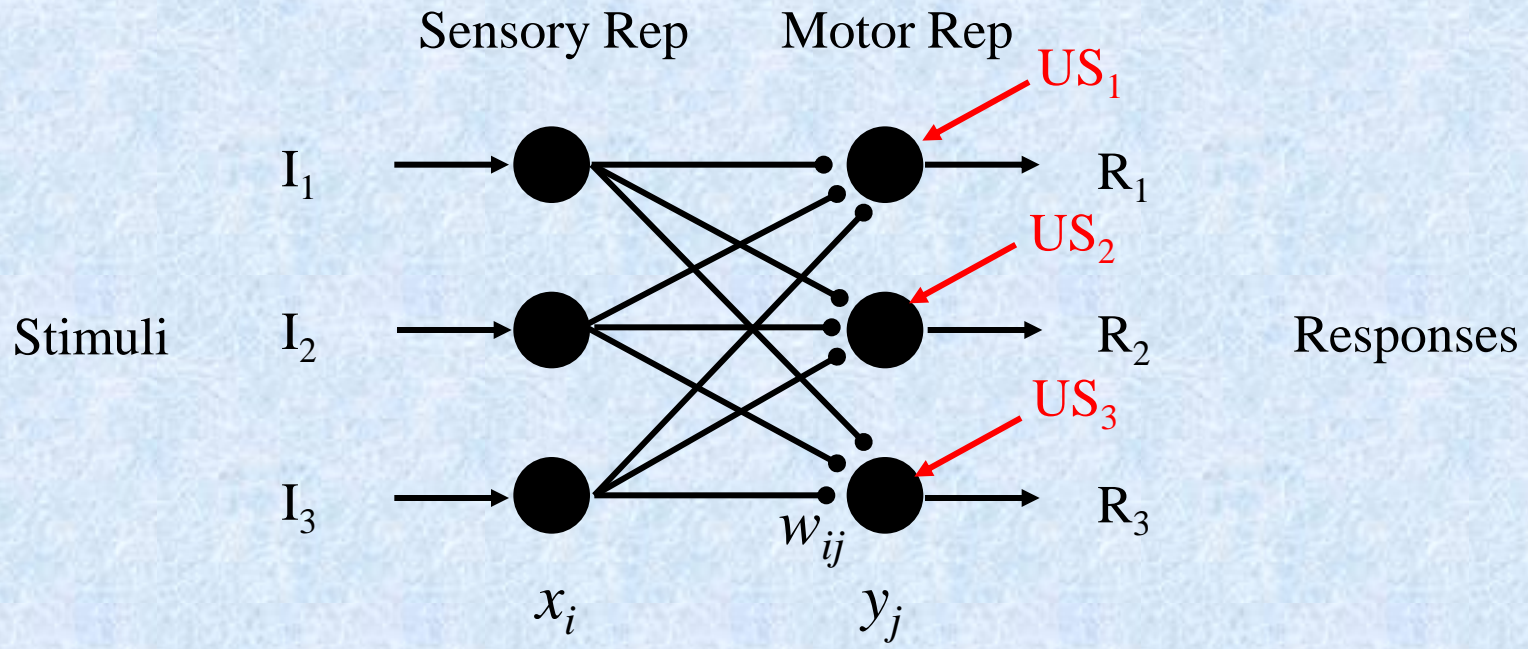


## (2) Memory interval: irrelevant events

Conditional probability  $p(\text{CR}|\text{CS})$  changes as a function of practice – nonstationary process

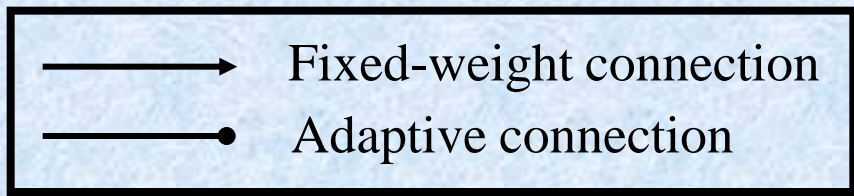
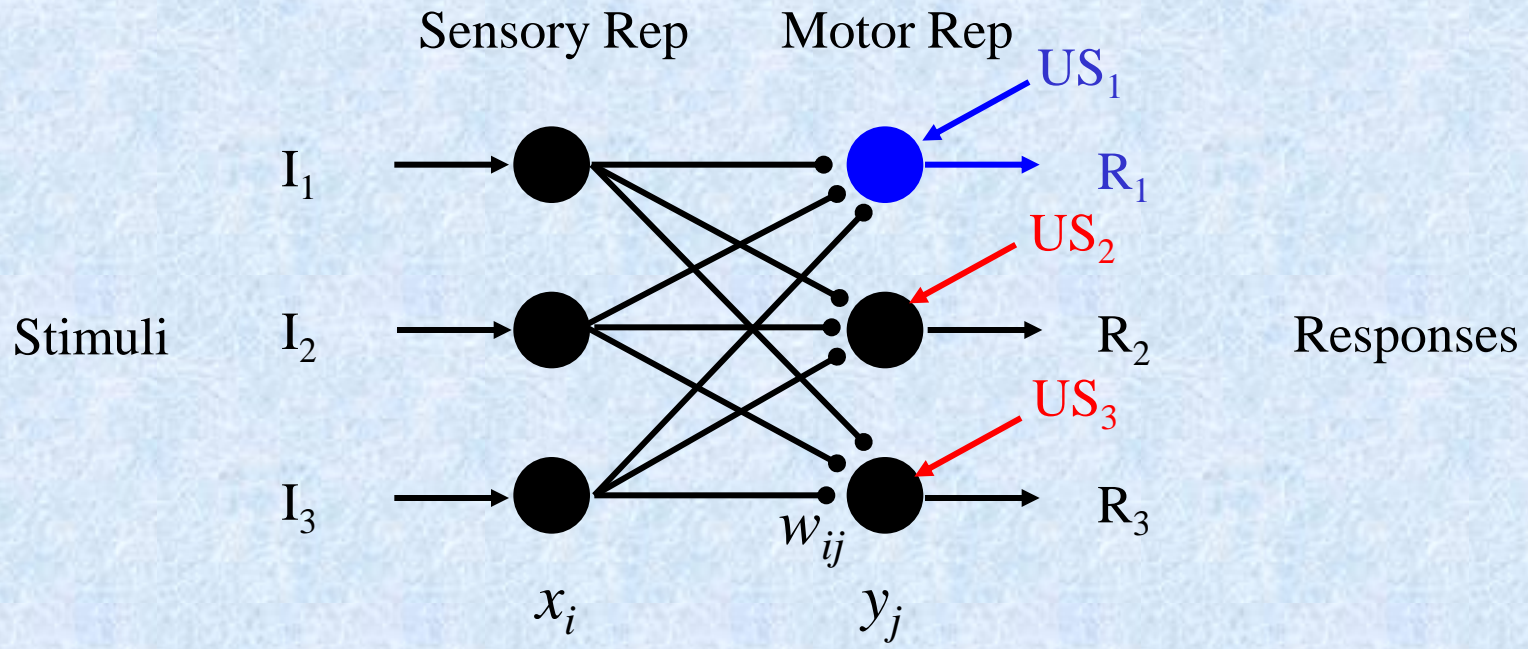
Change persists over long interval: LTM

Consider the following simple network, which includes a sensory representation and a motor representation:

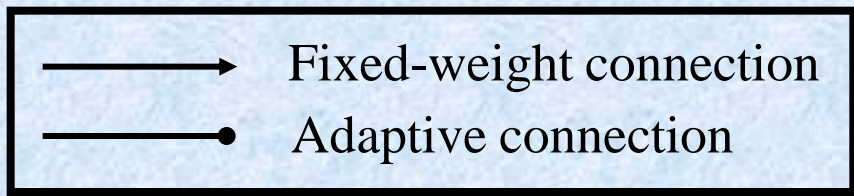
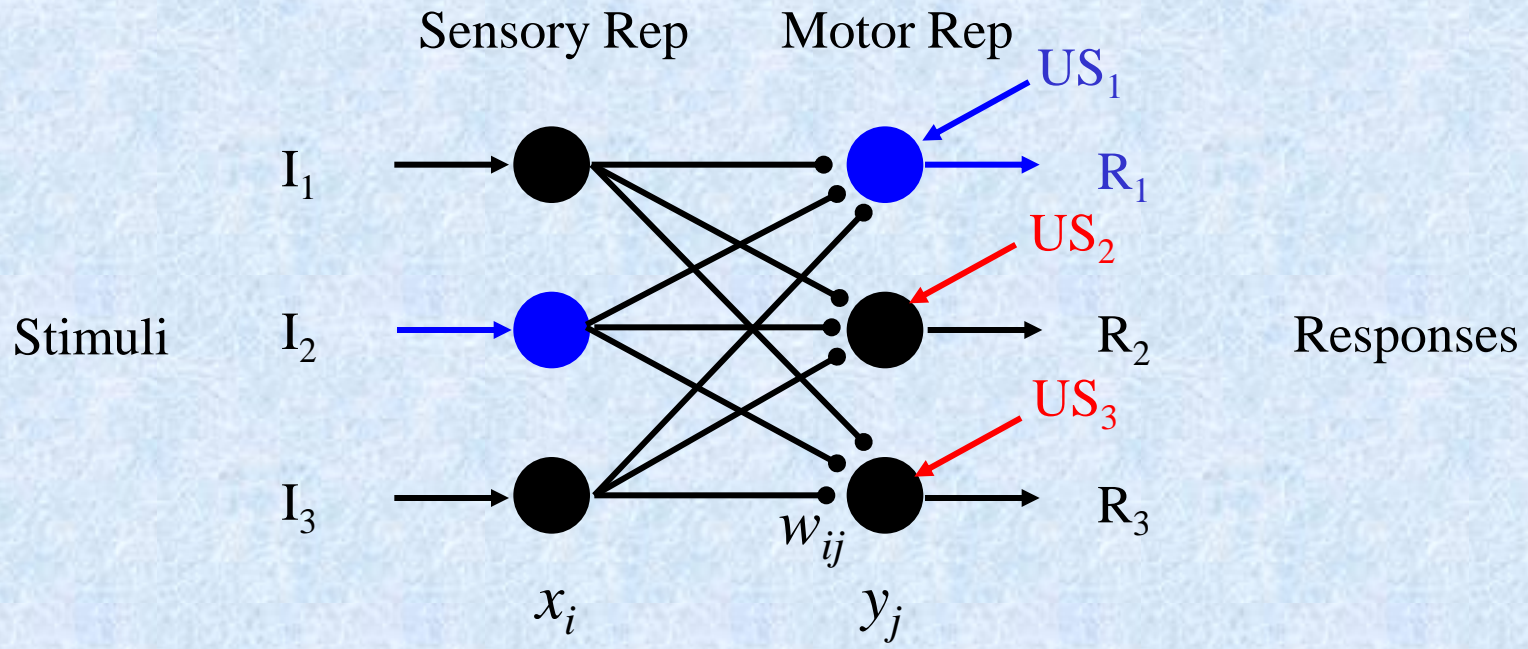




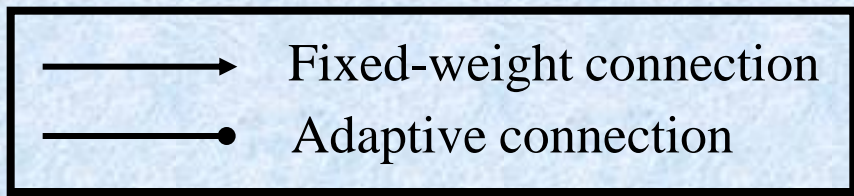
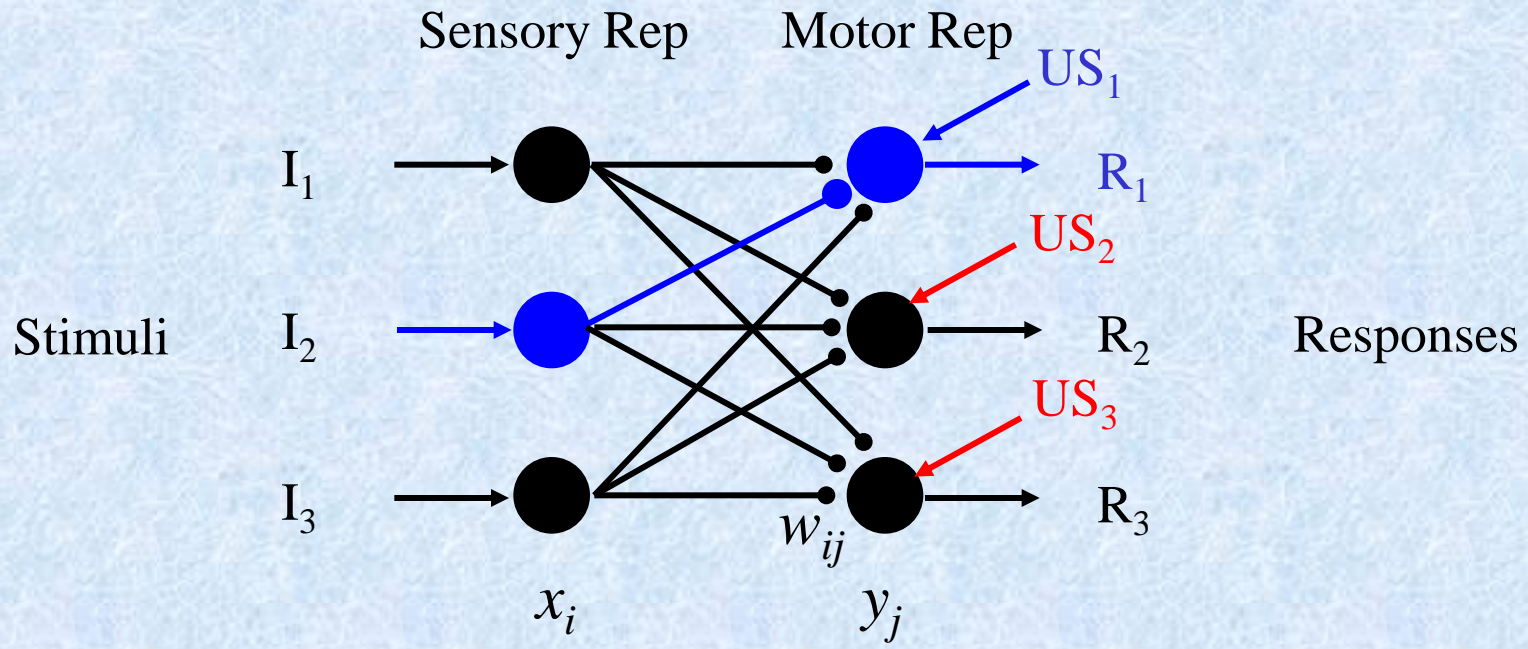
Consider the following simple network, which includes a sensory representation and a motor representation:



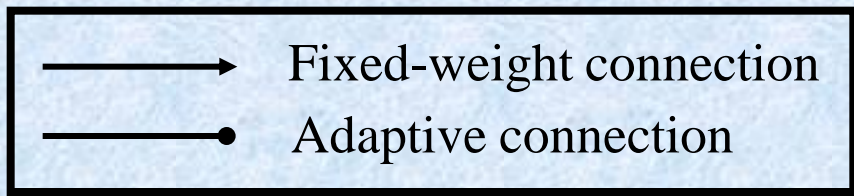
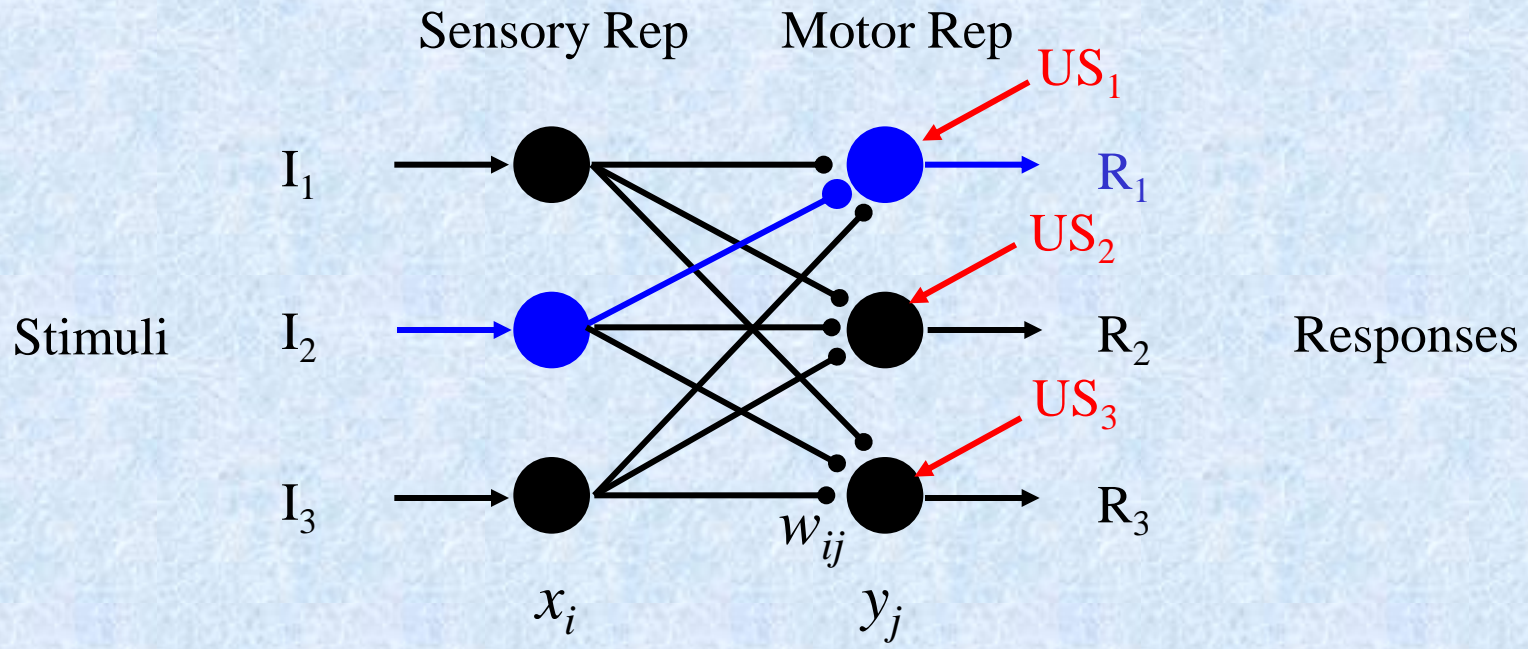
Consider the following simple network, which includes a sensory representation and a motor representation:



Consider the following simple network, which includes a sensory representation and a motor representation:



Consider the following simple network, which includes a sensory representation and a motor representation:





# Observations

A CS can be conditioned to an event that occurs after the stimulus has disappeared

A US that occurs shortly after a CS typically conditions more strongly and quickly than a US that occurs much later than the CS

If the US occurs at the same time or before the CS, then no conditioning takes place

What does this imply about the sensory representation?

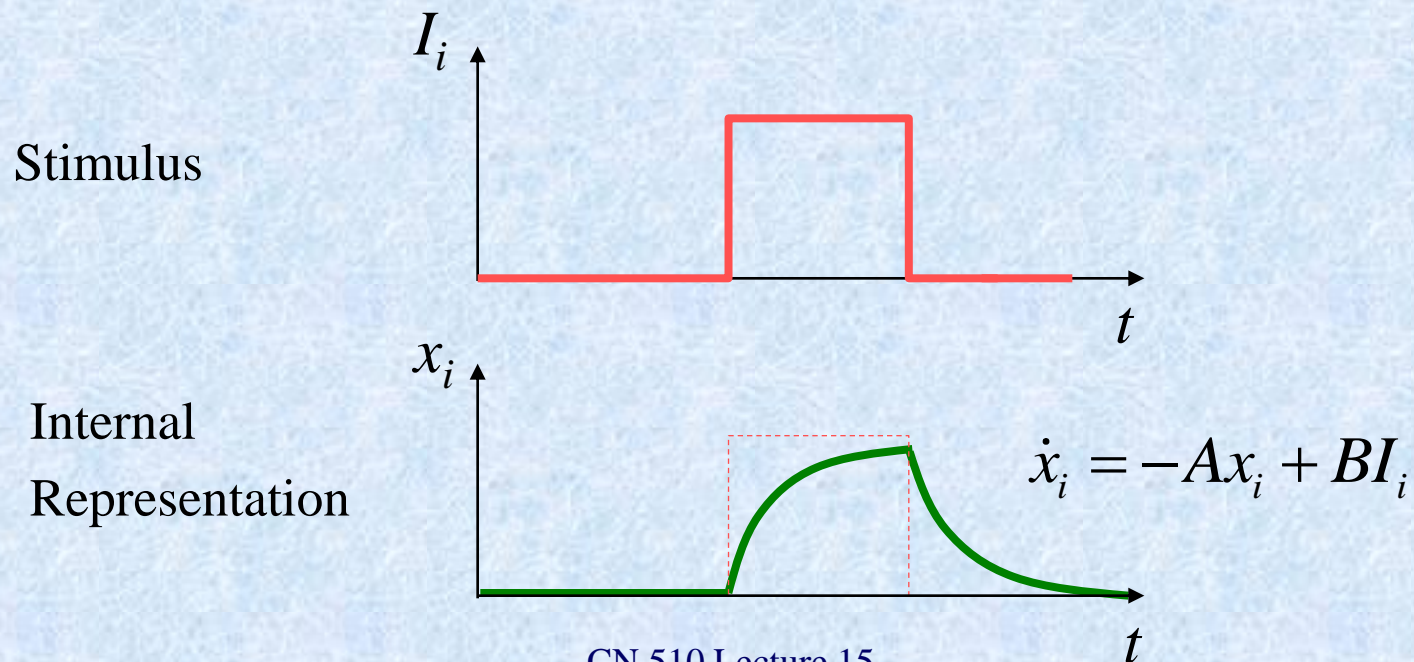
- The representation apparently persists after the stimulus is gone
- Time scale of this is much longer than STDP scale

What behavioral factors might underlie these properties?

- A real-world stimulus that occurs immediately before an important event is a more reliable predictor of the event than one that occurs much earlier
- A real-world stimulus that occurs after, or at the same time as, an important event is not useful for predicting the event

What kind of sensory representation might satisfy the following three constraints:

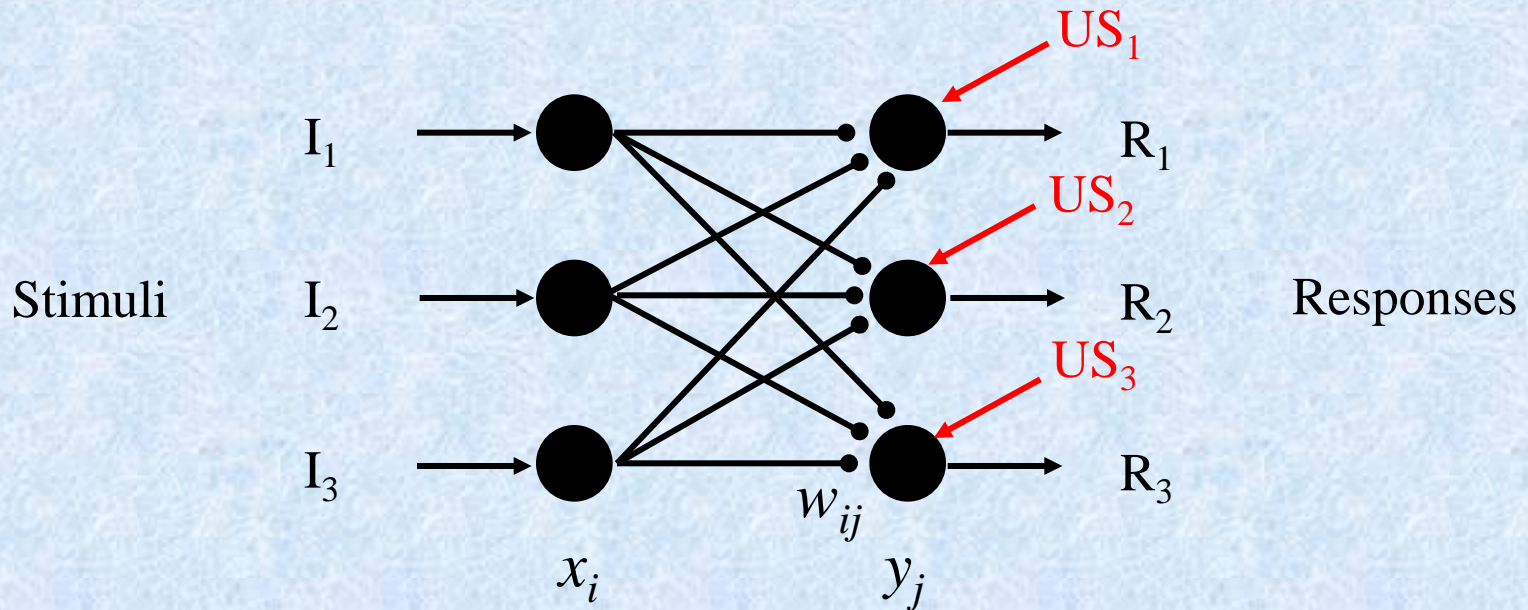
- the stimulus representation persists after the stimulus is gone,
- a US that occurs shortly after a CS typically conditions more strongly and quickly than a US that occurs much later than the CS,
- if the US starts at the same time or before the CS, then no conditioning takes place?



How long a time separation between CS and US can this network support?

Depends on the time constant of the leaky integrator, which is related to the parameter  $A$  (also called the decay rate) in the equation:

$$\dot{x}_i = -Ax_i + BI_i$$



$$\dot{x}_i = -Ax_i + BI_i$$

$$\dot{x}_i = -Ax_i + (B - x_i)I_i$$

The analogous equation for the motor representation:

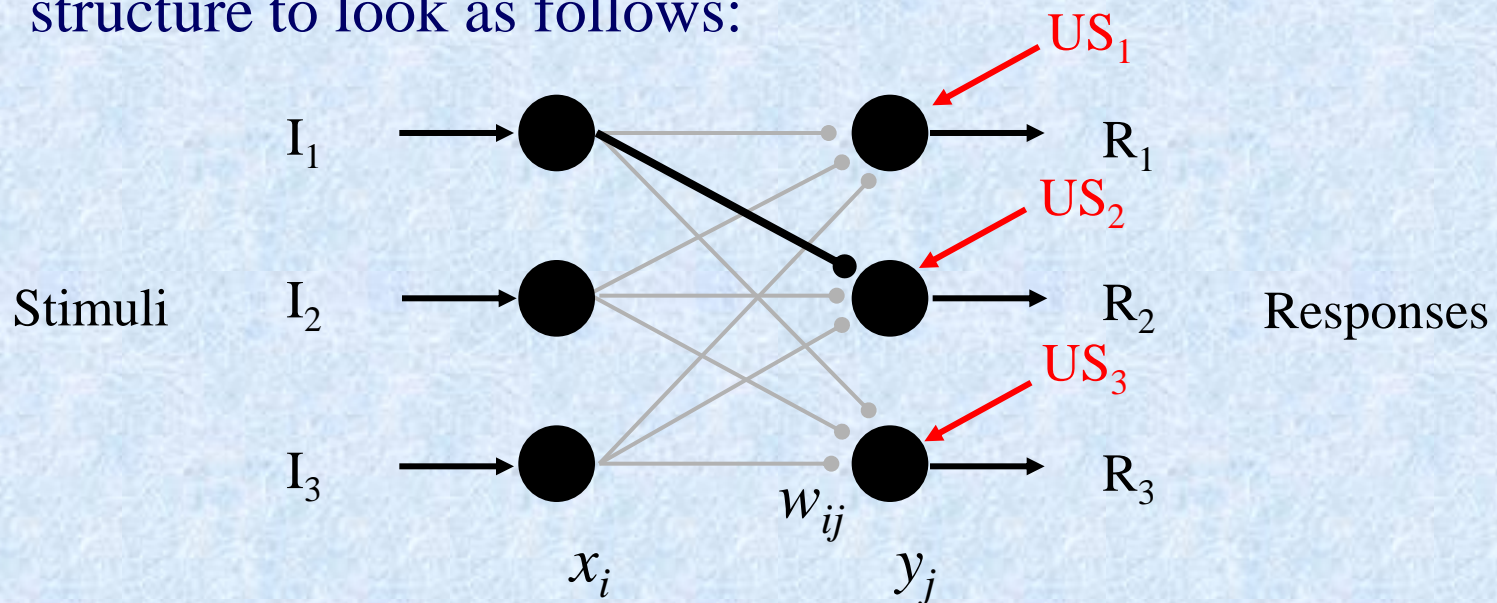
$$\dot{y}_j = -Ay_j + B \sum_i x_i w_{ij} + US_j$$

$$\dot{y}_j = -Ay_j + (B - y_j) \left( \sum_i x_i w_{ij} + US_j \right)$$



To explain classical conditioning, our network needs synaptic weights that encode learned associations between sensory stimuli and motor responses

If sensory stimulus  $I_1$  was repeatedly paired with a US that caused motor response  $R_2$ , but no other stimuli were consistently paired together, we would expect the network structure to look as follows:



What synaptic learning law would have this effect?

Passive decay:  $\dot{w}_{ij} = -Cw_{ij} + Dx_i y_j$

Association will die away if not continually practiced

However, in conditioning experiments, this does not typically happen: strong learned associations remain for long periods of time unless an animal undergoes extinction training:

CS: SHORT US: CS: SHORT US:  
BELL DELAY OMITTED ... BELL DELAY OMITTED ...

Here presynaptic signal is present and the association decays gradually, but without presynaptic signal it stays the same: sounds like presynaptically gated decay

# Presynaptically Gated Decay (a.k.a. “Outstar”)

$$\dot{w}_{ij} = -Cx_iw_{ij} + Dx_iy_j$$

Weights will decrease when a CS is encountered with no US

We can rewrite the equation as follows:

$$\dot{w}_{ij} = x_i(-Cw_{ij} + Dy_j)$$

Weight tracks the postsynaptic activity: for each CS it will increase for all paired US and decrease for all other US

# Postsynaptically Gated Decay (a.k.a. “Instar”)

Could we have used postsynaptically gated decay?

$$\dot{w}_{ij} = -C y_j w_{ij} + D x_i y_j$$

or

$$\dot{w}_{ij} = y_j (-C w_{ij} + D x_i)$$

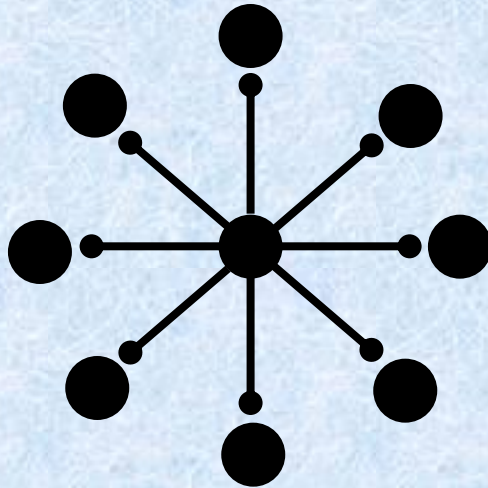
In this case, a response with no CS causes the weight to decay, while a CS with no response does not result in any change to the weight

This does not match the classical conditioning behavioral data as well as presynaptically gated decay, since responses in the absence of a CS (e.g., salivation in the case of Pavlov’s dogs) do not typically extinguish learned associations

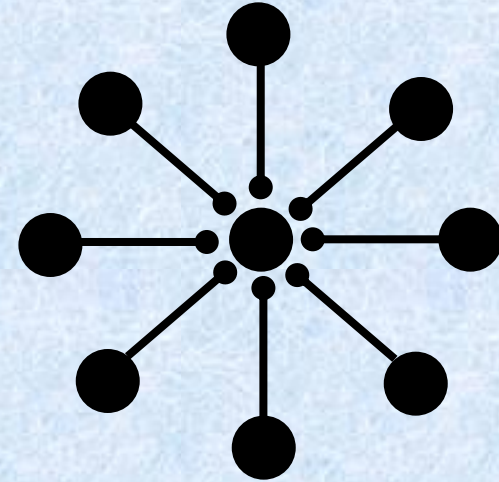


# Origin of the Terms “Instar” and “Outstar”

Outstar



Instar



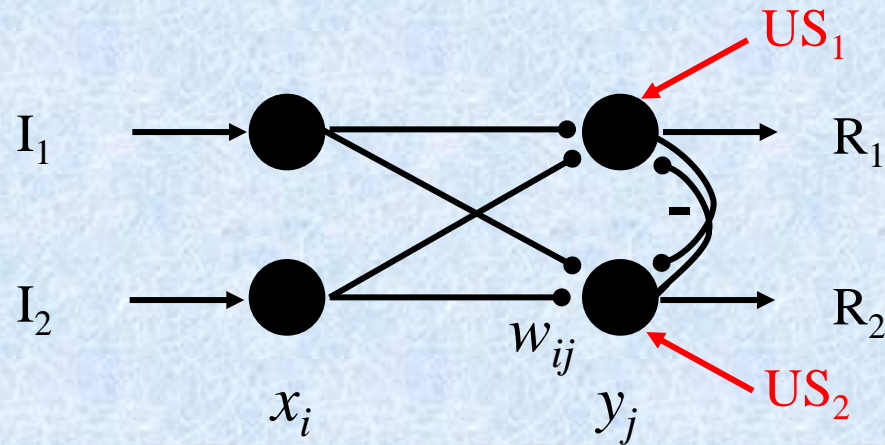
Weights projecting **from** a cell change when that cell is active – learns a pattern of outputs

Weights projecting **to** a cell change when that cell is active – learns to respond to a particular input pattern



# What Are the Units of Memory?

Is it sufficient to know weight change in a single pathway?

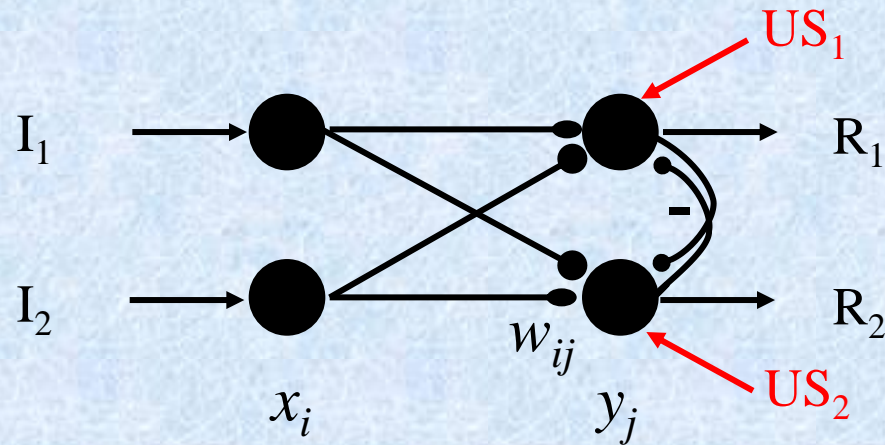


Let's say  $I_1$  co-occurs frequently with  $US_2$  and less frequently with  $US_1$

Let's say  $I_2$  co-occurs frequently with  $US_1$  and less frequently with  $US_2$

# What Are the Units of Memory?

Is it sufficient to know weight change in a single pathway?



All weights will grow with learning but  $w_{12}$  and  $w_{21}$  will grow more than  $w_{11}$  and  $w_{22}$

Functionally what is important is which input elicits which response, not which weight grew...

# Computing Normalized Weights

If we normalize weights

$\frac{W_{11}}{W_{12} + W_{11}}$  decreased and so did response

$\frac{W_{12}}{W_{12} + W_{11}}$  increased and so did response

So once again normalization seems like a solution

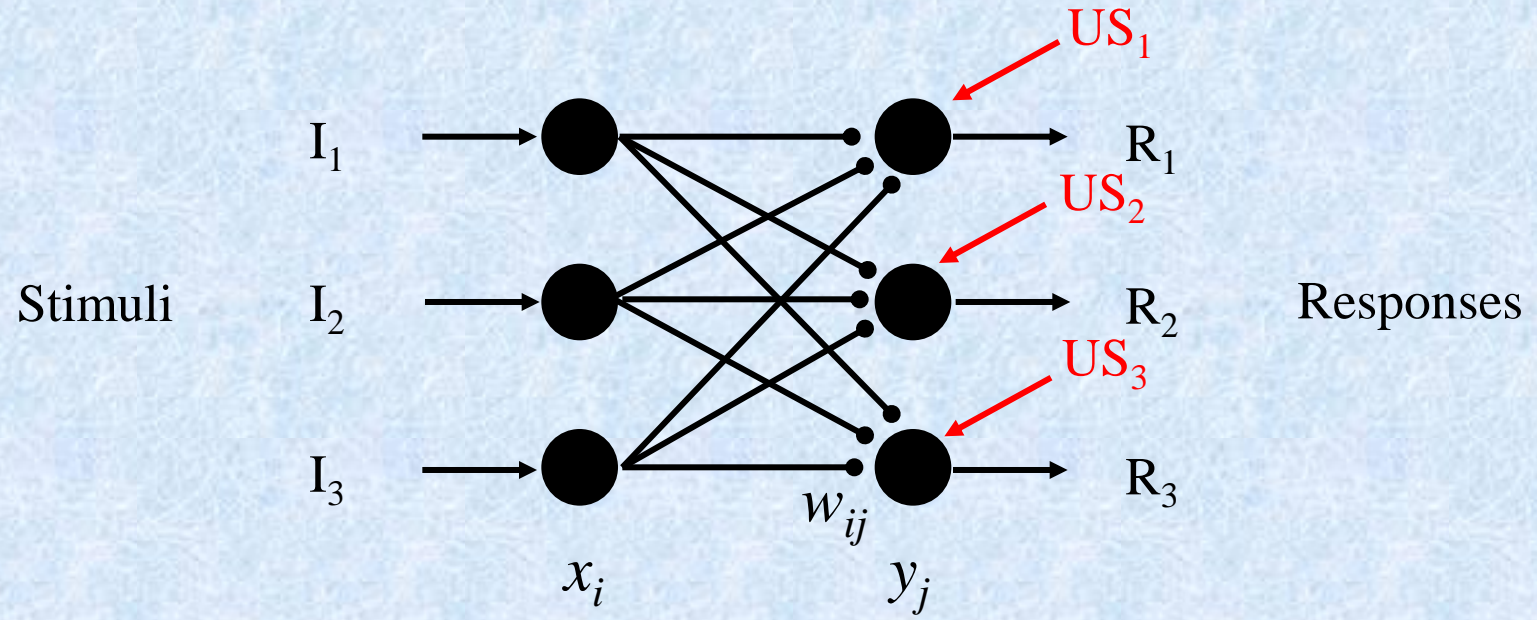
Bad news: functional significance is not local to a single synapse

Good news: knowing all outgoing weights of a cell gives an idea of what these weights have memorized

Also goes well with Levy's third rule of plasticity

Sometimes it's important to learn to suppress a response that leads to negative reinforcement, e.g. touching a hot stove

How could this be implemented in our network?



If a stimulus leads to a response that results in negative feedback, the synaptic weight between that stimulus' representation and the offending action should be weakened.

How could we modify our learning equation to handle this?

- Need some sort of **error signal** that changes what happens in the learning law when an error is produced:

$$\dot{w}_{ij} = -Dx_i w_{ij} + Ex_i y_j - F \cdot error$$

where  $F$  is another parameter. Is this sufficient?



With learning law on previous slide, all associations will be decreased when an error occurs, including correct associations unrelated to the current error.

$$\dot{w}_{ij} = -Dx_i w_{ij} + Ex_i y_j - F \cdot error$$

How can we fix this?

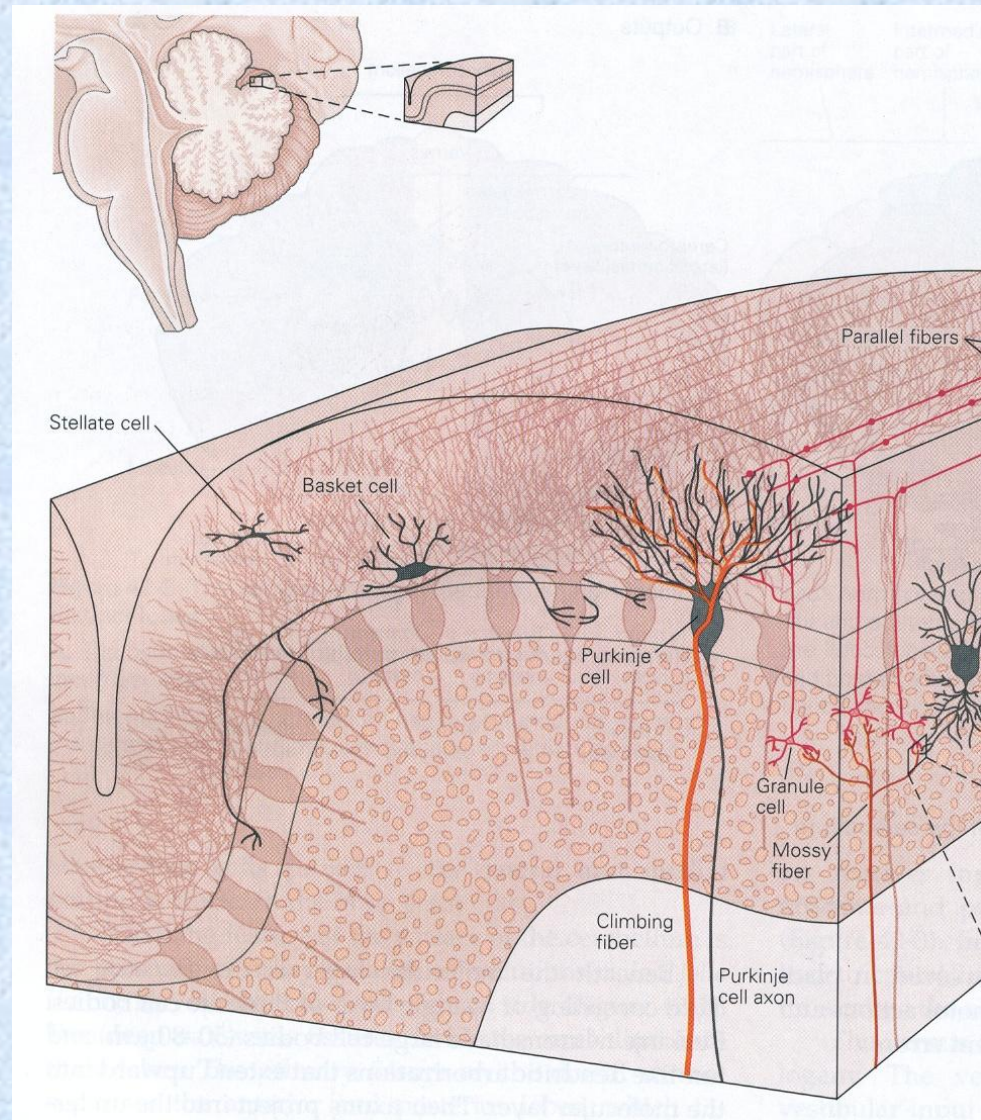
$$\dot{w}_{ij} = -Dx_i w_{ij} + Ex_i y_j - Fx_i y_j \cdot error$$

What condition on the parameters has to hold to account for the observation that negative reinforcement on a single trial can often permanently suppress an errant response in future trials?

A similar process appears to take place in the cerebellar cortex during motor learning:

- Errant action results in climbing fiber activation
- Climbing fiber activation causes “complex spike” in Purkinje cell that decreases the strengths of synapses from active parallel fibers (LTD)

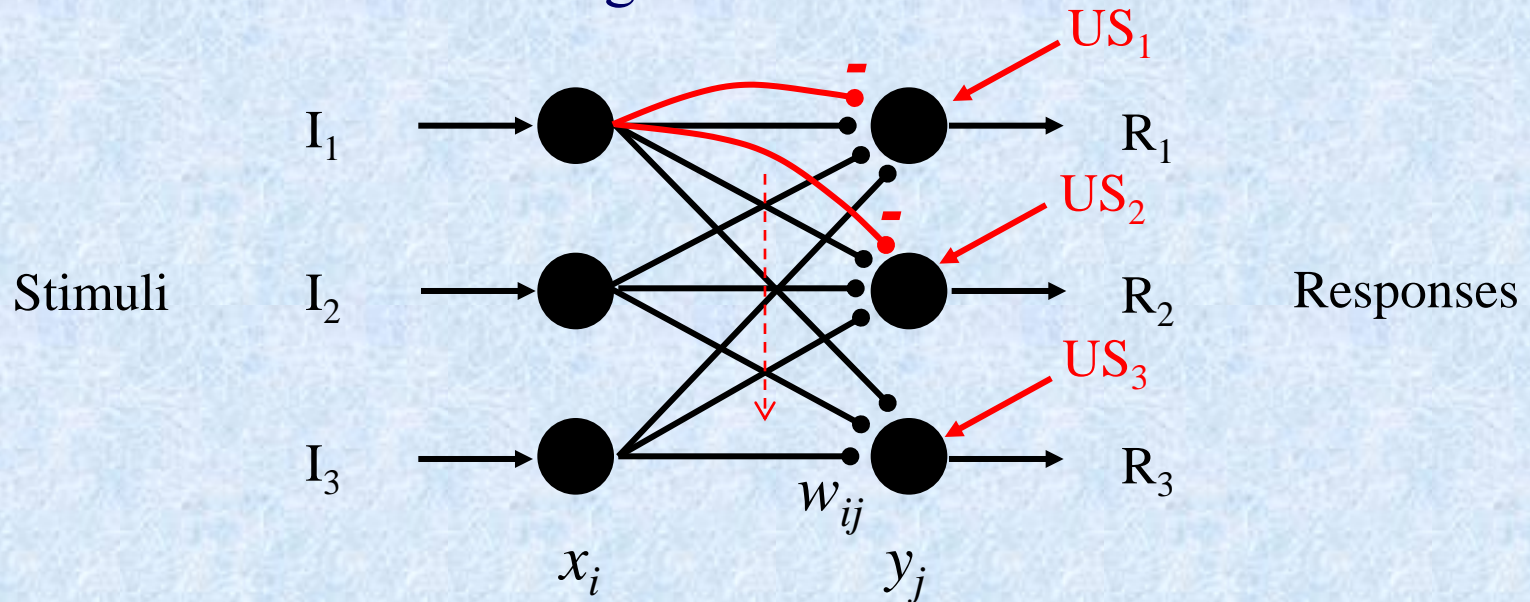
This learning is not as fast as the “single trial” learning discussed on the preceding slide (smaller value of  $F$  in the equation).





# Alternative Solution

Could alternatively have added inhibitory connections to our network to encode “negative” associations:



$$\dot{y}_j = -Ay_j + B \sum_i x_i w_{ij}^+ - F \sum_k x_k w_{kj}^-$$

$$\dot{y}_j = -Ay_j + (B - y_j) \sum_i x_i w_{ij}^+ - (F + y_j) \sum_k x_k w_{kj}^-$$

In this case, we do not need the error term in the equation for the excitatory connections:

$$\dot{w}_{ij}^+ = -Dx_i w_{ij}^+ + Ex_i y_j$$

What equation should be used to produce weight increases for the inhibitory connections?

$$\dot{w}_{ij}^- = Hx_i y_j \cdot error$$

If we added a decay term to this equation, what form should it take?

$$\dot{w}_{ij}^- = -Gx_i y_j w_{ij}^- + Hx_i y_j \cdot error$$

# Model's Weaknesses

Sensory stimuli are “lumped”; i.e., one cell per stimulus

Motor stimuli are “lumped”

The model might not produce a good quantitative fit to behaviorally measured learning curves, e.g.

- Time constant of leaky integrator most likely is not sufficient to account for “inverted U” learning curve

Limited comparison of model cells to those in a real biological nervous system

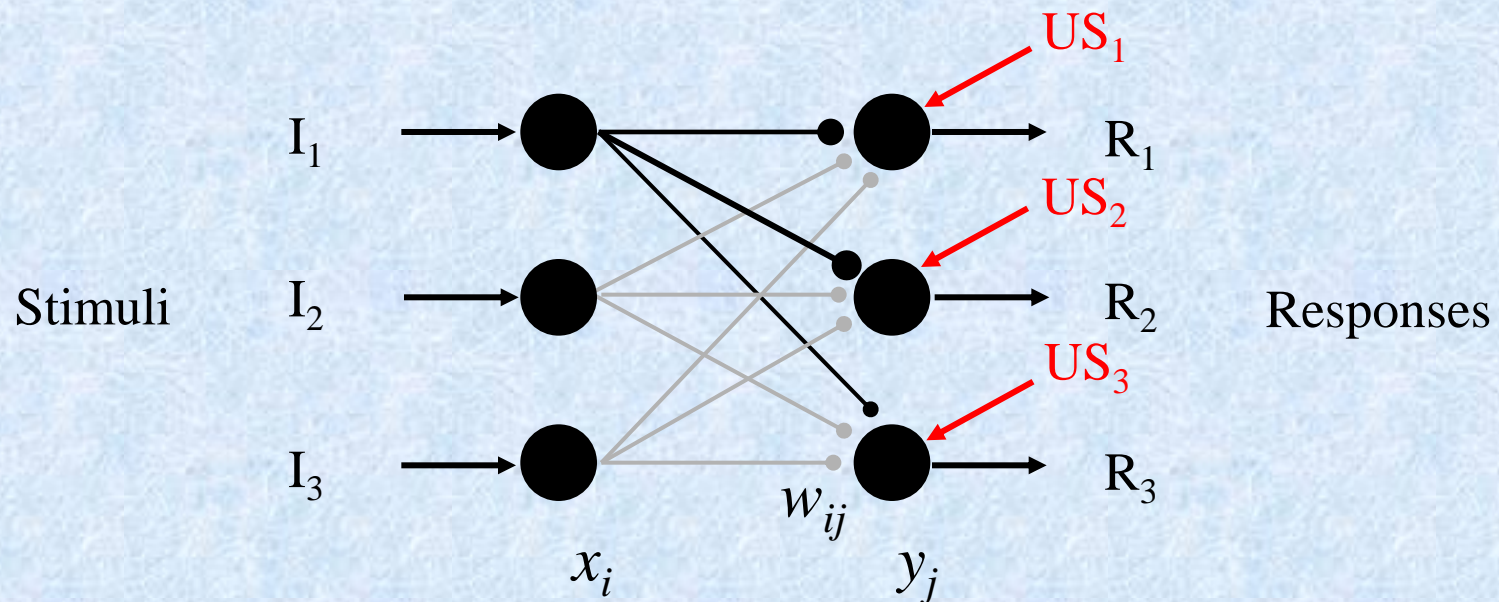
Please note, that this model totally disregards reinforcement part of learning, reward value, or predictive error



# Generalizing

It is often useful to adjust the weights so that they reflect a pattern of activity (given that this pattern is set by external input)

Veridical prediction or recall will happen if the relative sizes of weights  $w_{ij}$  resemble the relative sizes or frequencies of the sampled  $y_j$

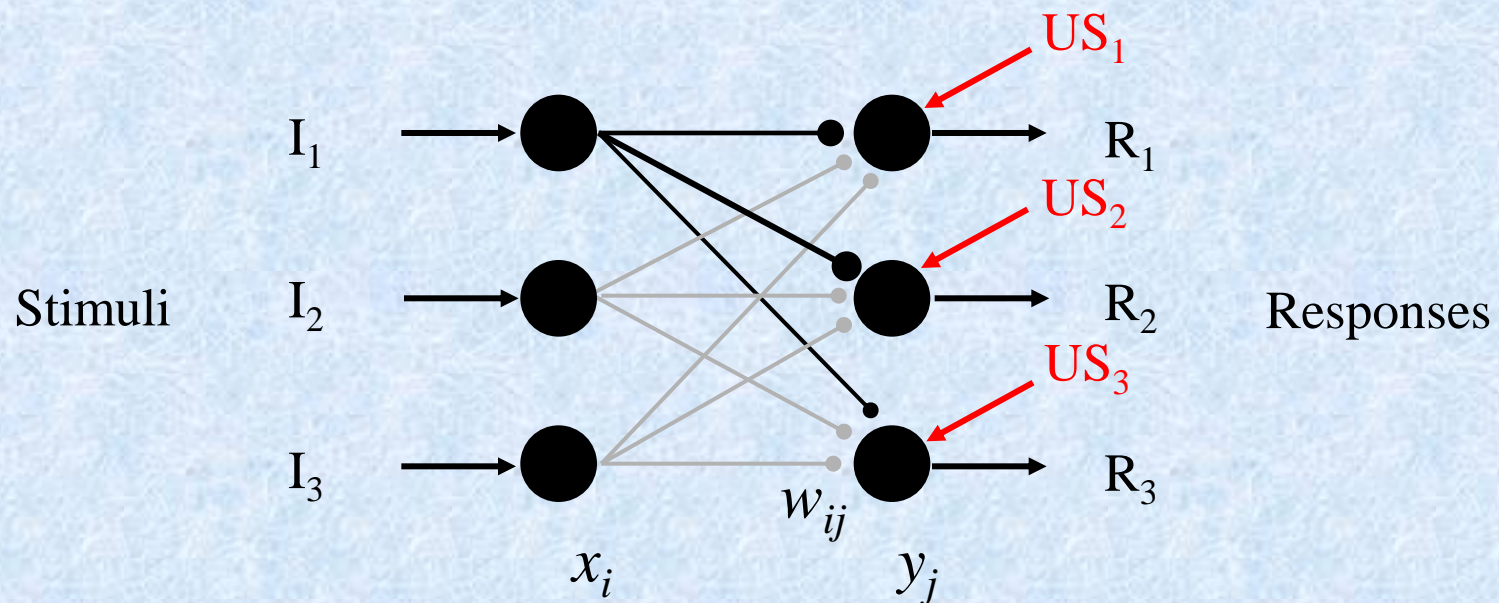


# Generalizing

Note that the proper operation of this network consists of two phases:

Learning (Encoding) – when synapses learn, but do not transmit the information

Recall (Retrieval) – when synapses transmit but are not modified

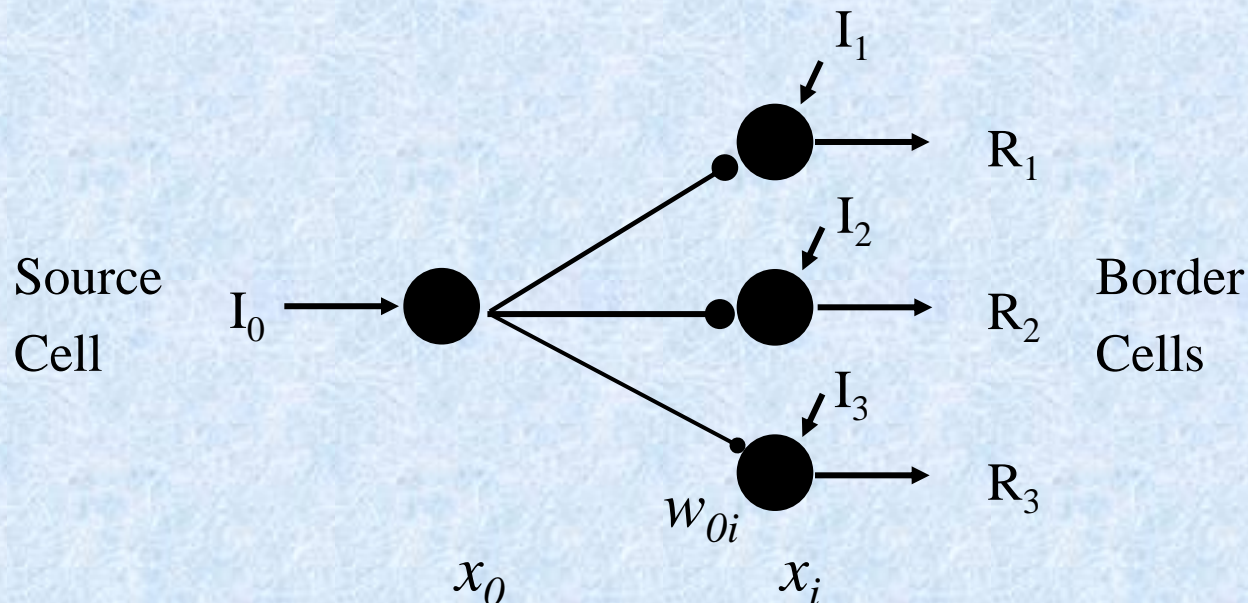


# Outstar Network

The source cell samples the border and learns the spatial pattern across it

Learning depends on the sampling regime (activation regime of the source cell) because we use presynaptically gated rule

$$\dot{w}_{0i} = x_0 (-Cw_{0i} + Dx_i)$$

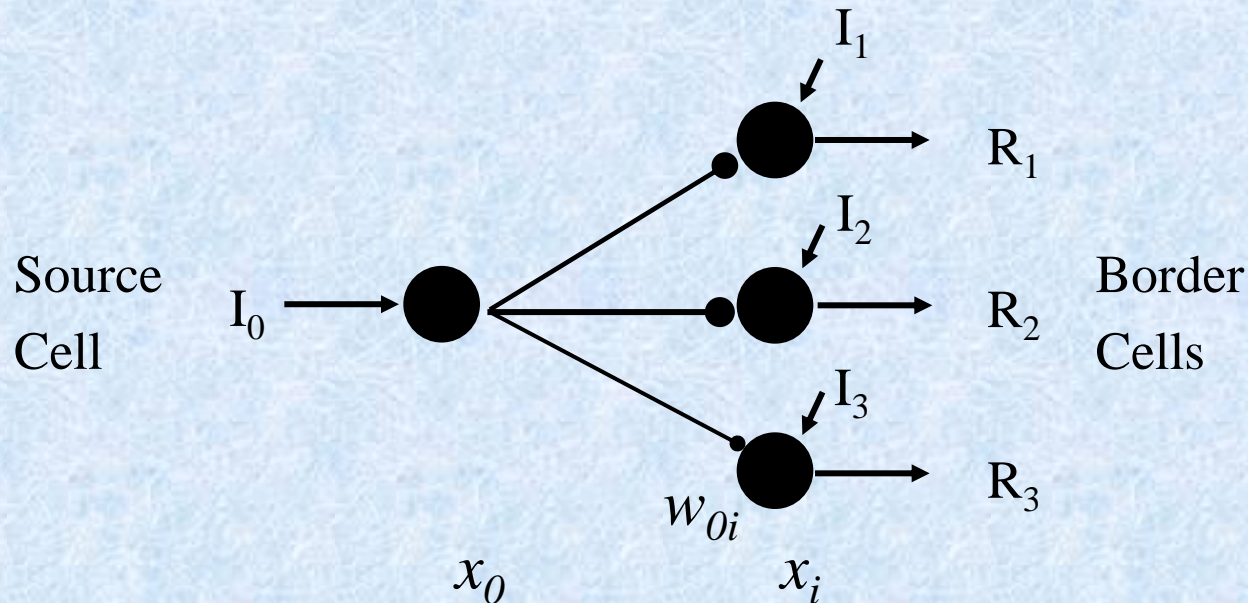


# Generalized Outstar Network

$$\dot{x}_0 = -a_0 x_0 + I_0(t)$$

$$\dot{x}_i = -a_i x_i + b_{0i} [x_0 - \Gamma_{0i}]_+ w_{0i} + I_i(t)$$

$$\dot{w}_{0i} = [x_0 - \Gamma_{0i}]_+ (-c_{0i} w_{0i} + d_{0i} x_i)$$



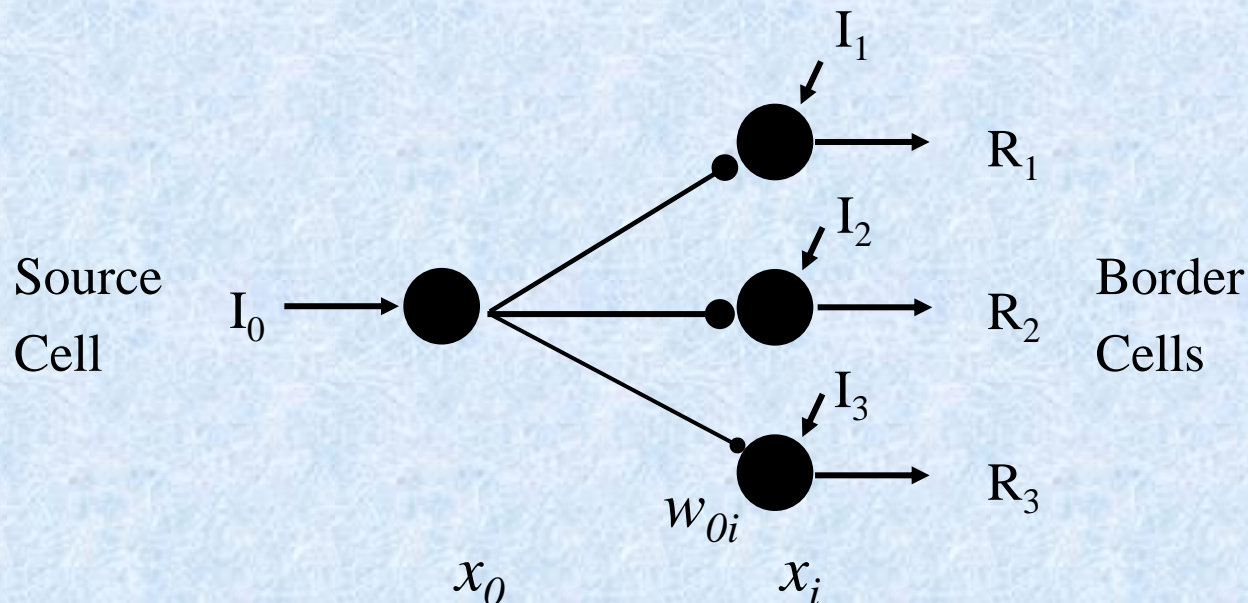
# Unbiased Outstar Network

Unbiased outstar has all non-adaptive parameters the same for all border cells:

$$\dot{x}_0 = -a_0 x_0 + I_0(t)$$

$$\dot{x}_i = -a x_i + b [x_0 - \Gamma]_+ w_{0i} + I_i(t)$$

$$\dot{w}_{0i} = [x_0 - \Gamma]_+ (-c w_{0i} + d x_i)$$





# What Can Outstar Learn?

One sampling with very strong learning rate: eidetic memory

Some samplings of similar patterns with moderate learning rate: pattern average or prototype

Constant sampling: time-average of input (useless?)

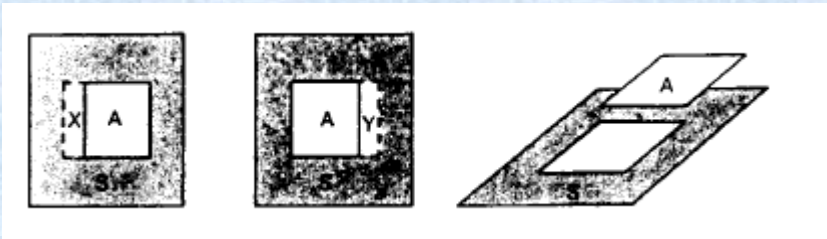
Does eidetic memory exist *in vivo*?

Bela Julesz (pronounced you-less) claimed that it does in some rare subjects...

# Random Dot Stereogram Method



Disparity information specifying figure-in-depth exists only when two (high-dimensional) spatial patterns, one appropriate to each eye, simultaneously perturb the visual system



Day one:

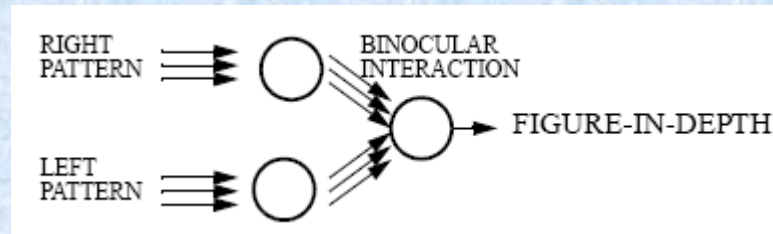
- Present right random-dot image to right eye of subject

Day two:

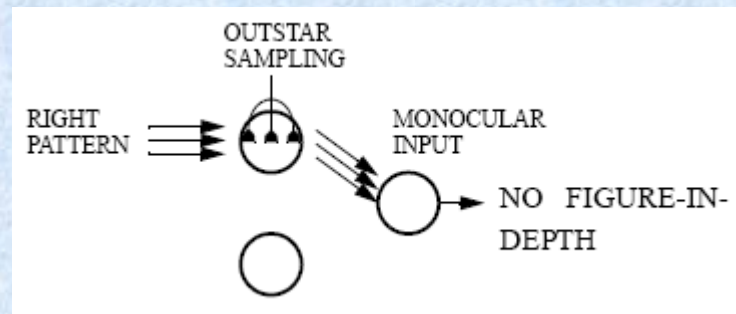
- Present left random-dot image to left eye of subject
- Request recall of right image

# Outstar Explanation

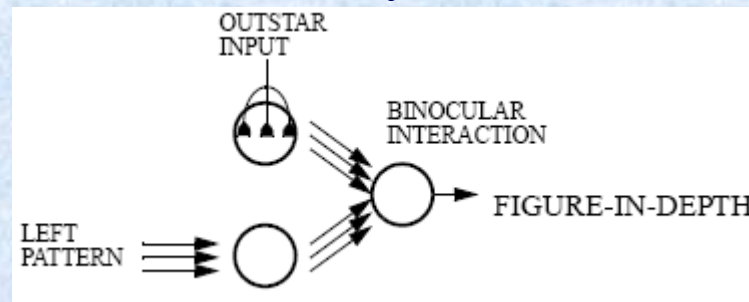
Normal use of stereogram:



Day 1:



Day 2:



# Outstar Learning Theorem

Unbiased outstar with  $b$ ,  $c$ , and  $d$  subsuming  $[x_0 - \Gamma]_+$ :

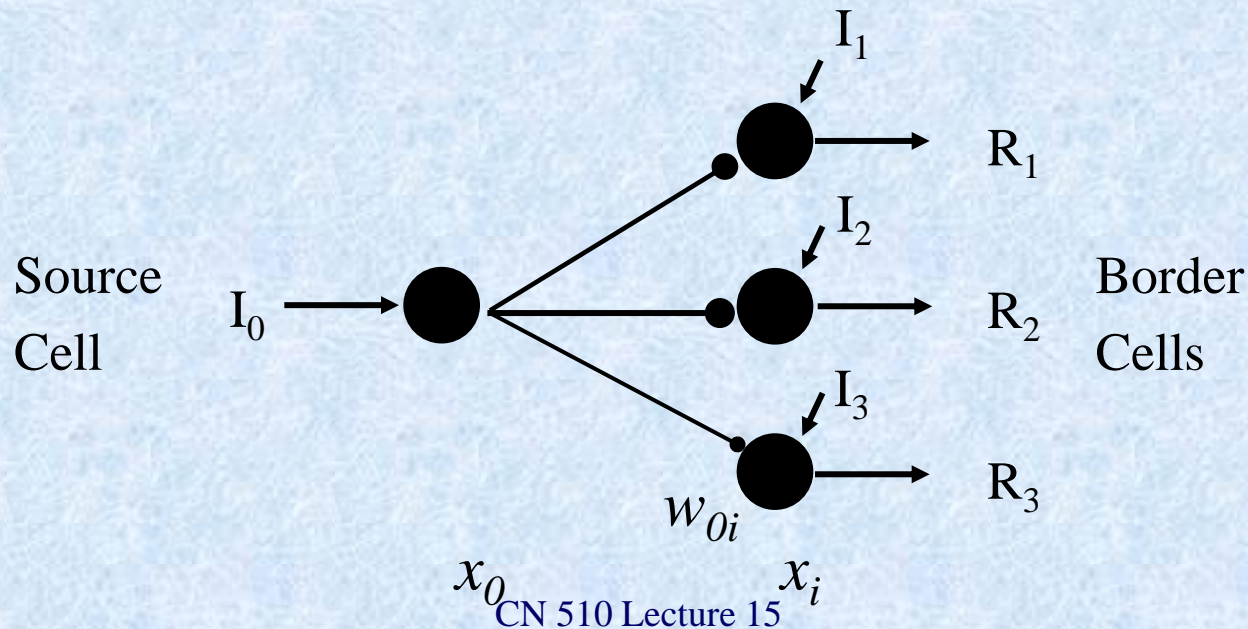
$$\dot{x}_0 = -a_0 x_0 + I_0(t)$$

$$\dot{x}_i = -a(t)x_i + b(t)w_{0i} + I_i(t)$$

$$\dot{w}_{0i} = -c(t)w_{0i} + d(t)x_i$$

For generality  $a$  is also a function of time

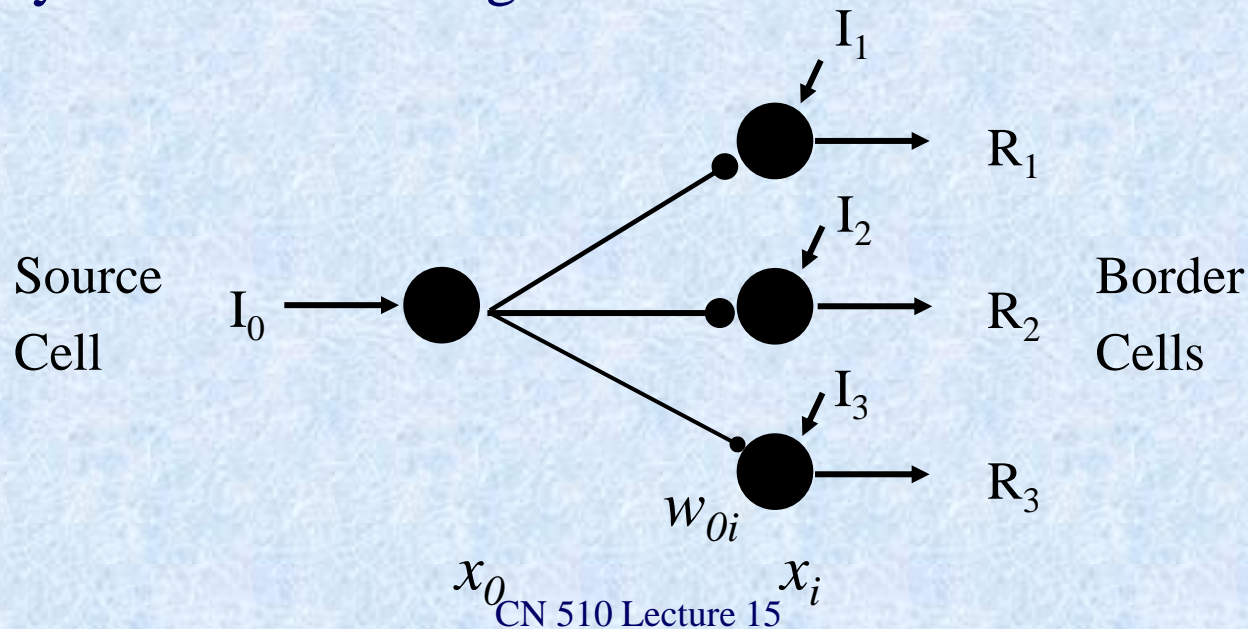
$$a, b, c, d \geq 0$$



# Outstar Learning Theorem

$$\begin{aligned}\dot{x}_0 &= -a_0 x_0 + I_0(t) \\ \dot{x}_i &= -a(t) x_i + b(t) w_{0i} + I_i(t) \\ \dot{w}_{0i} &= -c(t) w_{0i} + d(t) x_i\end{aligned}$$

Note that  $c(t)$  can represent not only presynaptically gated decay, but also passive decay or postsynaptically gated decay without violating the theorem



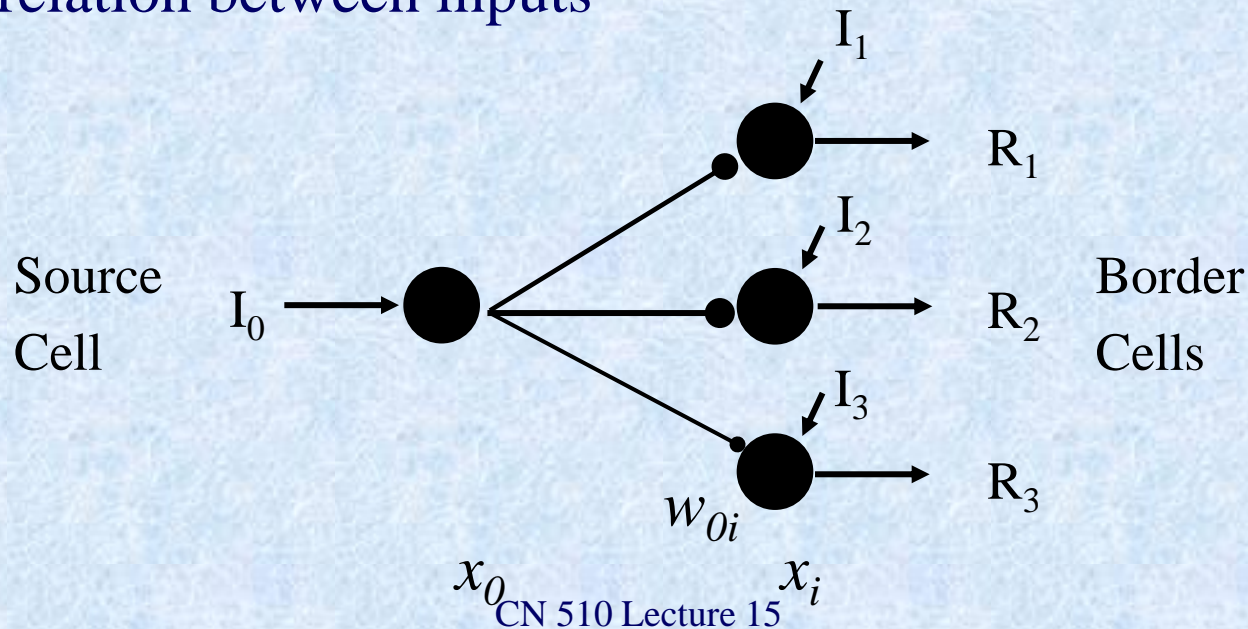


# Outstar Learning Theorem

Specify a measure  $\Theta$  that captures the spatial pattern of the input to be learned

Specify a measure  $W$  that captures the spatial pattern of the weights that encode the input

Prove that for appropriate sampling regime  $W$  approaches  $\Theta$  regardless of the initial activation, weight distribution, or correlation between inputs



# Pattern Variables

How large is the input at site  $i$  relative to the total input

$$\Theta_i = \frac{I_i}{\sum_{k=1}^n I_k}$$

How large is the weight at site  $i$  relative to all weights

$$W_{0i} = \frac{w_{0i}}{\sum_{k=1}^n w_{0k}}$$

Prove that

$$\lim_{t \rightarrow \infty} W_{0i} = \Theta_i$$

In practice, inputs can only affect weight values indirectly, via the activations

Because of this, we also need to define a measure  $X$ , the spatial pattern of activations,

$$X_i = \frac{x_i}{\sum_{k=1}^n x_k}$$

and show that:

$$\lim_{t \rightarrow \infty} W_{0i} = \lim_{t \rightarrow \infty} X_i = \Theta_i$$

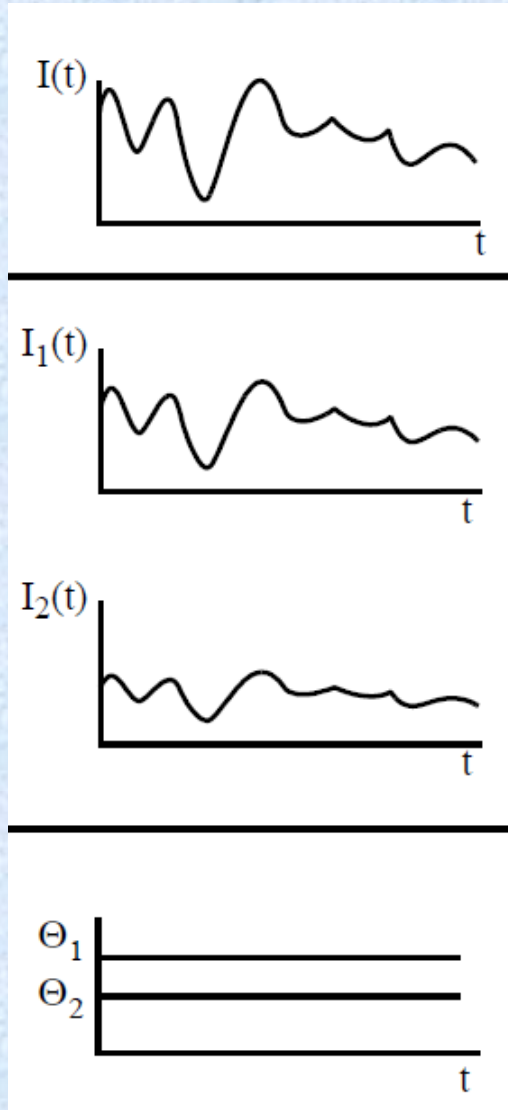
or

$$\lim_{t \rightarrow \infty} \vec{W} = \lim_{t \rightarrow \infty} \vec{X} = \vec{\Theta}$$

# Pattern vs Energy Variables

Note that since the pattern variables are normalized by definition

$$\sum_{i=1}^n W_{0i} = \sum_{i=1}^n X_i = \sum_{i=1}^n \Theta_i = 1$$



# Assumptions

Assume that activations and weights are bounded from above by a certain number  $M$

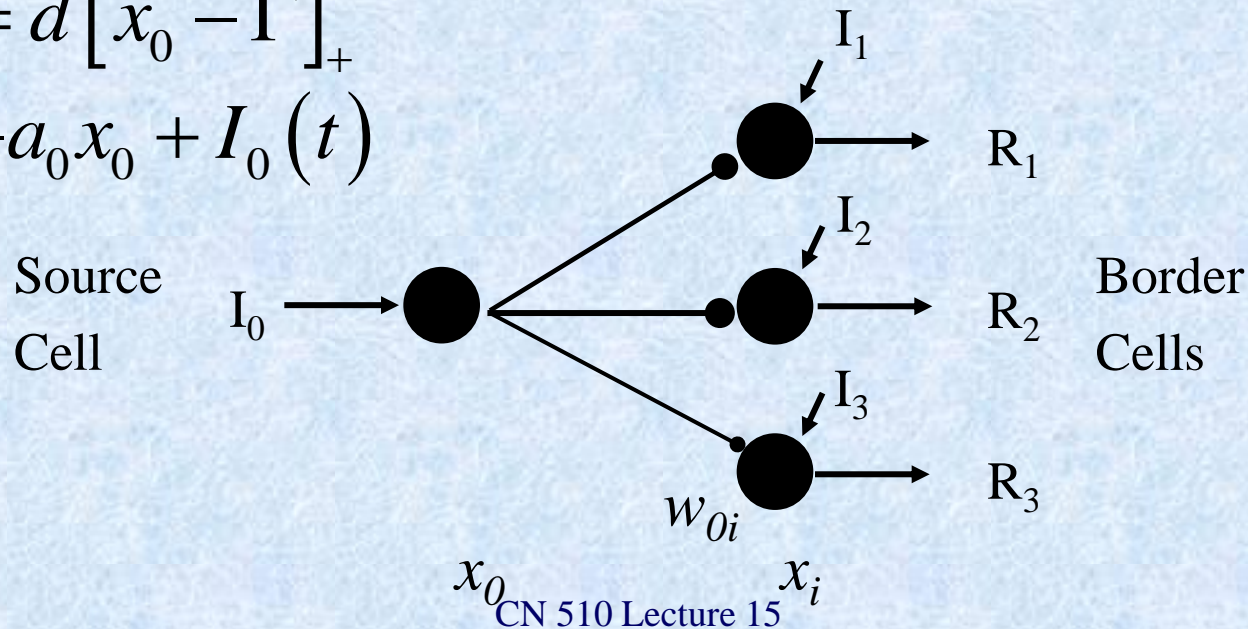
Assume that CS ( $I_0$ ) is presented sufficiently often to give enough time to sample the input

$$\int_0^{\infty} d(t) dt = \infty$$

$$\dot{w}_{0i} = -c(t) w_{0i} + d(t) x_i$$

$$d(t) = d[x_0 - \Gamma]_+$$

$$\dot{x}_0 = -a_0 x_0 + I_0(t)$$

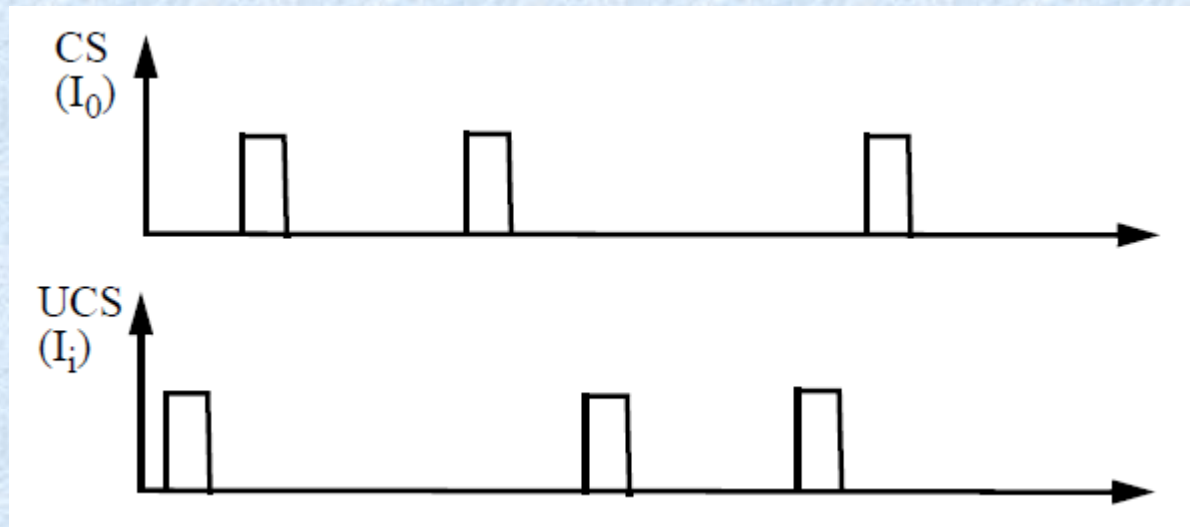




# Assumptions

Input pattern ( $I_i$ ) is presented sufficiently often so that the decay rate is slow enough for activations to be above zero at all times

This allows for uncorrelated CS and US presentations



# Six Step Intuitive Proof

1. Verify positivity of  $x$  and  $w$  to make sure pattern variables are valid
2. Derive equations for total energy variables
3. Derive a system in terms of pattern variables
4. Double check what happens if there is no  $I_0$
5. Double check what happens if there is no  $I_i$
6. Analyze pattern variable system under normal conditions

# Positivity of Solutions

To justify use of the pattern variables

$$X_i = \frac{x_i}{\sum_{k=1}^n x_k} \quad W_{0i} = \frac{w_{0i}}{\sum_{k=1}^n w_{0k}}$$

we must ensure that  $x$  and  $w$  are positive

Assume they are at  $t=0$

In

$$\dot{x}_i = -a(t)x_i + b(t)w_{0i} + I_i(t)$$

$b(t)$  and  $I(t)$  are non-negative

If we assume for a moment that  $w$  is non-negative, then

$$\dot{x}_i \geq -a(t)x_i$$

# Positivity of Solutions

$$\dot{x}_i \geq -a(t)x_i$$

In integral form this is

$$x_i \geq x_i(0)e^{\int_0^t -a(v)dv} > 0$$

Thus, given non-negative  $w$ , if  $x$  starts positive it stays positive

But the same argument is true for

$$\dot{w}_{0i} = -c(t)w_{0i} + d(t)x_i$$

Thus for positive initial conditions on  $x$  and  $w$  these variables stay positive



# Total Energy Equations

Sum

$$\dot{x}_i = -a(t)x_i + b(t)w_{0i} + I_i(t)$$

across all  $x_i$

$$\sum_{i=1}^n \dot{x}_i = -a \sum_{i=1}^n x_i + b \sum_{i=1}^n w_{0i} + \sum_{i=1}^n I_i$$

or

$$\dot{x} = -ax + bw + I$$

Similarly,

$$\dot{w}_{0i} = -c(t)w_{0i} + d(t)x_i$$

becomes

$$\dot{w} = -cw + dx$$

# System in Pattern Variables

From

$$X_i = \frac{x_i}{\sum_{k=1}^n x_k} = \frac{x_i}{x}$$

we can compute

$$\frac{dX_i}{dt} = \frac{d}{dt} \left( \frac{x_i}{x} \right) = \frac{1}{x^2} (x\dot{x}_i - x_i\dot{x}) = \frac{1}{x} \left( \dot{x}_i - \frac{x_i}{x} \dot{x} \right)$$

and substitute

$$\dot{x}_i = -ax_i + bw_{0i} + I_i$$

$$\dot{x} = -ax + bw + I$$

$$\frac{dX_i}{dt} = \frac{1}{x} \left( -ax_i + bw_{0i} + I_i - \frac{x_i}{x} (-ax + bw + I) \right)$$

$$\frac{dX_i}{dt} = \frac{1}{x} \left( -\mathbf{ax}_i + bw_{0i} + I_i + \mathbf{ax}_i - \frac{1}{x} (bx_iw + Ix_i) \right)$$

$$\frac{dX_i}{dt} = \frac{b}{x} \left( w_{0i} - \frac{x_iw}{x} \right) + \frac{1}{x} \left( I_i - \frac{Ix_i}{x} \right)$$

$$\frac{dX_i}{dt} = \frac{bw}{x} \left( \frac{w_{0i}}{w} - \frac{x_i}{x} \right) + \frac{I}{x} \left( \frac{I_i}{I} - \frac{x_i}{x} \right)$$

$$\frac{dX_i}{dt} = \frac{bw}{x} \left( \frac{w_{0i}}{w} - \frac{x_i}{x} \right) + \frac{I}{x} \left( \frac{I_i}{I} - \frac{x_i}{x} \right)$$

$$\frac{dX_i}{dt} = \frac{bw}{x} (W_{0i} - X_i) + \frac{I}{x} (\Theta_i - X_i)$$

Finally, let  $A = \frac{bw}{x}$  and  $B = \frac{I}{x}$

$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$



Similarly,

$$W_{0i} = \frac{w_{0i}}{w} \quad \frac{dW_{0i}}{dt} = \frac{1}{w} \left( \dot{w}_{0i} - \frac{w_{0i} \dot{w}}{w} \right)$$

$$\frac{dW_{0i}}{dt} = \frac{1}{w} \left( -Cw_{0i} + dx_i + Cw_{0i} - dw_{0i}x \frac{1}{w} \right)$$

$$\frac{dW_{0i}}{dt} = \frac{dx}{w} \left( \frac{x_i}{x} - \frac{w_{0i}}{w} \right) = \frac{dx}{w} (X_i - W_{0i})$$

and using  $C = \frac{dx}{w}$  we get  $\dot{W}_{0i} = C(X_i - W_{0i})$

## System in Pattern Variables

$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

$$A = \frac{bw}{x} \geq 0 \quad B = \frac{I}{x} \geq 0 \quad C = \frac{dx}{w} \geq 0$$

Note that the decay rates of the original system do not affect the system in pattern variables

Energies as well as parameters of the original system do not affect the direction of pattern variable changes, they only affect the rates of these changes

## System without $I_0$

$$\dot{W}_{0i} = C(X_i - W_{0i}) \quad C = \frac{d(t)x}{w} \quad d(t) = d[x_0 - \Gamma]_+$$

In the absence of the sampling signal  $x_0$  goes to 0, so  $d(t)$  goes to 0, so  $C$  goes to 0

$$\dot{W}_{0i} = 0$$

So there is a perfect memory of the spatial pattern that we have learned until this point

Note that this does not mean that the weights themselves stay constant; this will depend on the decay gating in the original system

## System without US

$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

$$A = \frac{bw}{x} \qquad B = \frac{I}{x} \qquad C = \frac{dx}{w}$$

If  $I=0$  then  $B=0$  and the system reduces to

$$\dot{X}_i = A(W_{0i} - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

and if we have learned the pattern ( $W_{0i} = X_i = \Theta_i$ ), then we will keep it after the input shuts off



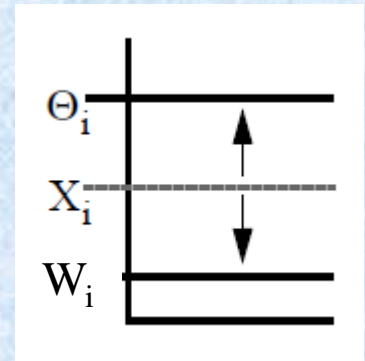
# Convergence of the System

$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

Intuitively the system converges due to the transitivity of tracking

But what if  $X$  hangs between  $W$  and  $\Theta$ ?

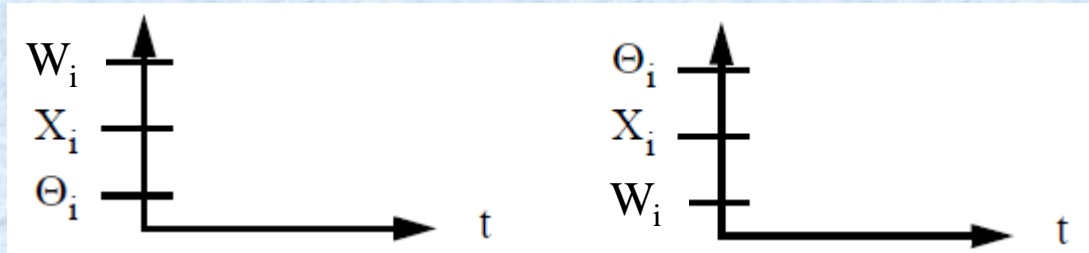


Should not be a problem as  $W$  tracks  $X$ , and thus its ability to pull  $X$  away from  $\Theta$  is continually eroding

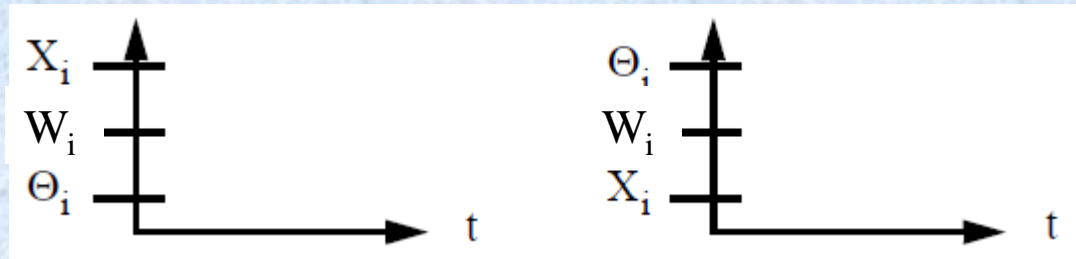
Let's make sure it is so

# Three Cases

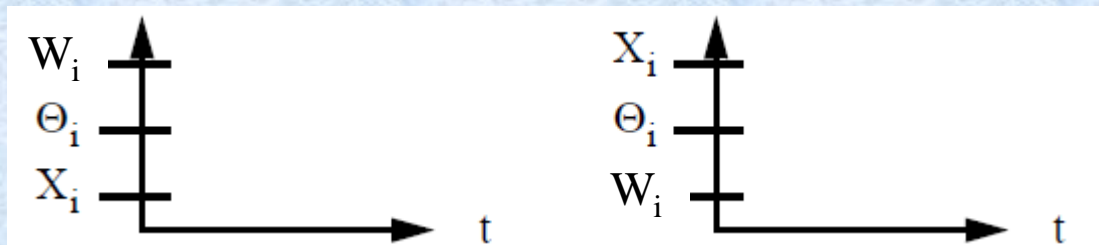
$X$  and  $W$  on the same side of  $\Theta$  with  $X$  closer to  $\Theta$



$X$  and  $W$  on the same side of  $\Theta$  with  $W$  closer to  $\Theta$

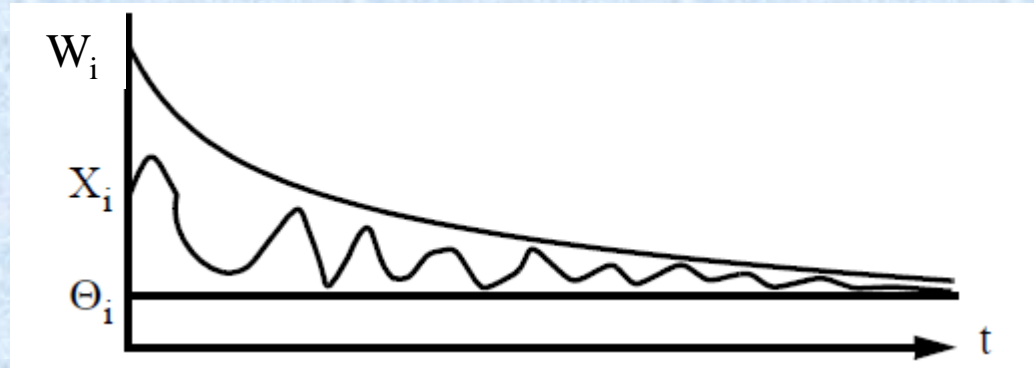


$X$  and  $W$  on different sides of  $\Theta$



# Case 1

$X$  and  $W$  on the same side of  $\Theta$  with  $X$  closer to  $\Theta$



$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

$X$  cannot go above  $W$  or below  $\Theta$ , unless  $W$  crosses  $\Theta$

$W$  cannot cross  $\Theta$  unless  $X$  does

So  $X$  has to stay between the two, and  $W$  is monotonically decreasing

$$\lim_{t \rightarrow \infty} \vec{W} = \lim_{t \rightarrow \infty} \vec{X} = \vec{\Theta}$$

## Case 2

$X$  and  $W$  on the same side of  $\Theta$  with  $W$  closer to  $\Theta$



$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

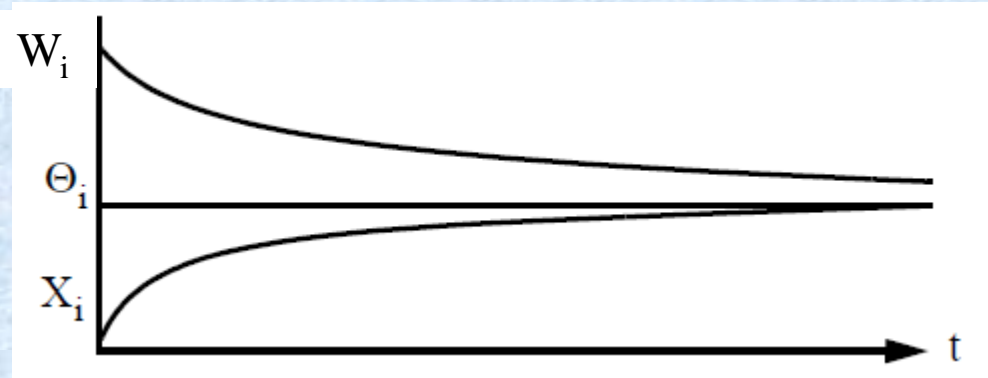
$X$  is increasing, and  $W$  is decreasing, until they cross, and then we are back to case 1

$$\lim_{t \rightarrow \infty} \vec{W} = \lim_{t \rightarrow \infty} \vec{X} = \vec{\Theta}$$



## Case 3

$X$  and  $W$  on different sides of  $\Theta$



$$\dot{X}_i = A(W_{0i} - X_i) + B(\Theta_i - X_i)$$

$$\dot{W}_{0i} = C(X_i - W_{0i})$$

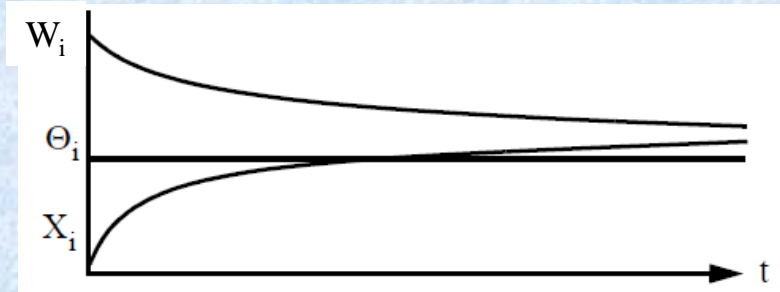
$X$  is increasing and  $W$  is decreasing

Eventually, one of the three things will happen

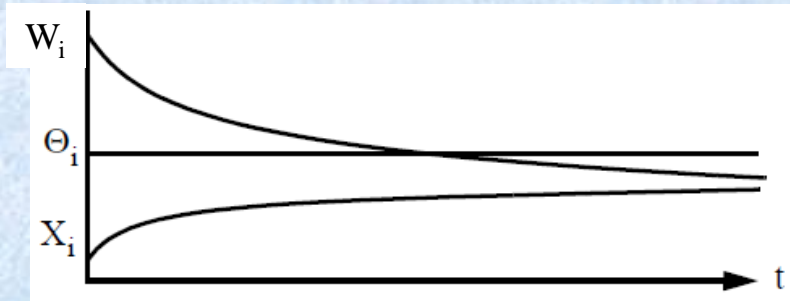
### Case 3

$$\lim_{t \rightarrow \infty} \vec{W} = \lim_{t \rightarrow \infty} \vec{X} = \vec{\Theta}$$

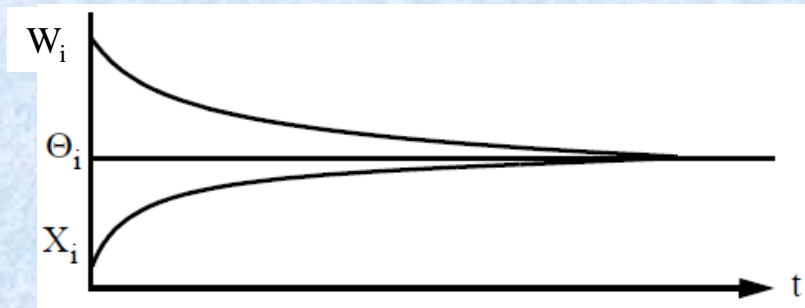
$X$  will cross  $\Theta$  and we are back to case 1



$W$  will cross  $\Theta$  and we are back to case 2



$X$  and  $W$  simultaneously reach  $\Theta$



# Summary

Outstar Learning Theorem proves that perfect learning is possible in terms of pattern variables so that

$$\lim_{t \rightarrow \infty} \vec{W} = \lim_{t \rightarrow \infty} \vec{X} = \vec{\Theta}$$

It requires

- no bias or omission in the input sampling (all  $b(t)$  are equal and non-zero)
- no other patterns at the border cells during sampling
- sufficient time to sample this particular pattern

These ideal conditions are hard to meet in practice, this is usually the case for most proven neural theorems

# Next Time

Chaining outstars for spatio-temporal learning  
Other models of sequence learning