

EECS 391

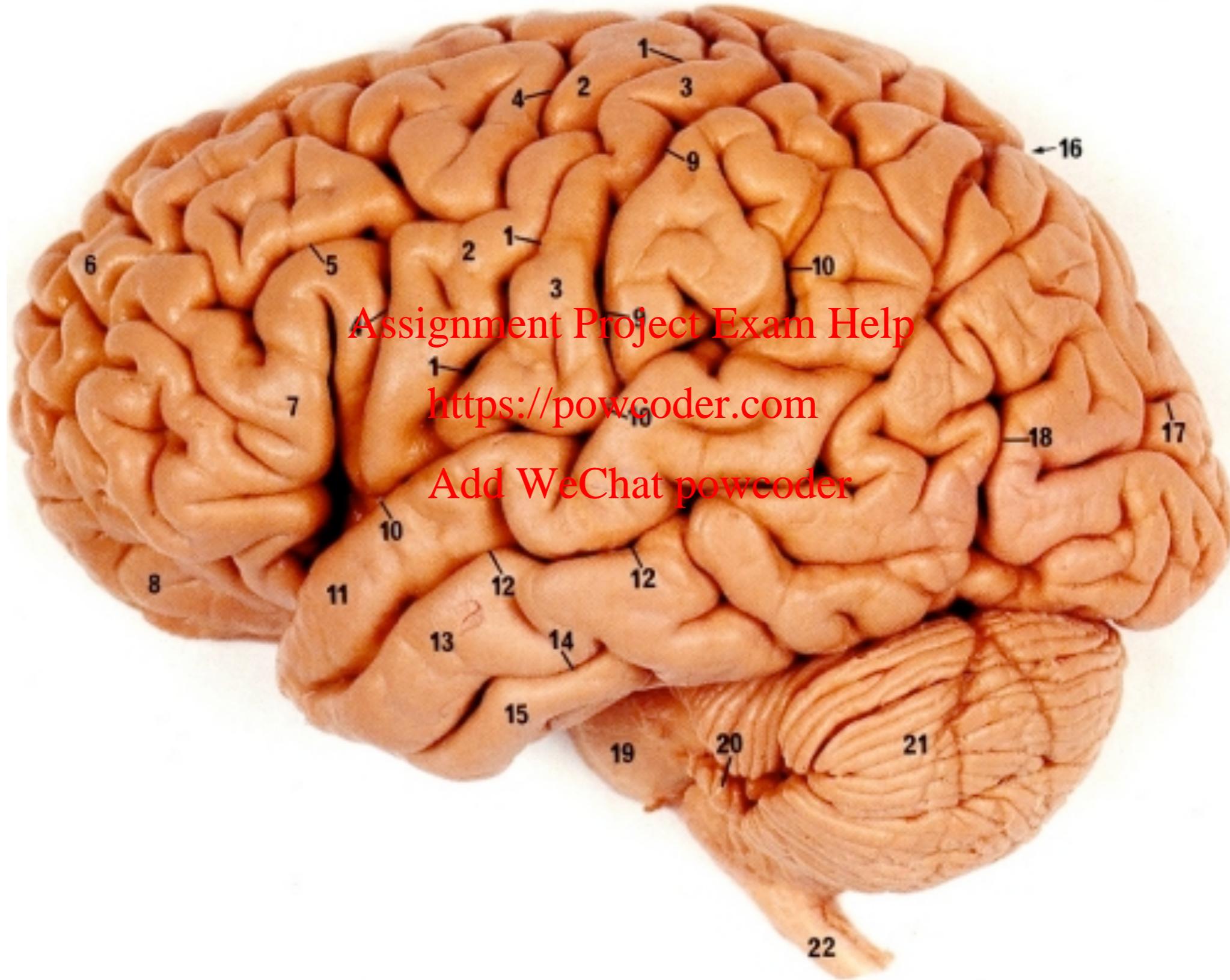
Intro to AI

Deep Belief Networks

<https://powcoder.com>

Add WeChat powcoder

L16 Thu Nov 2



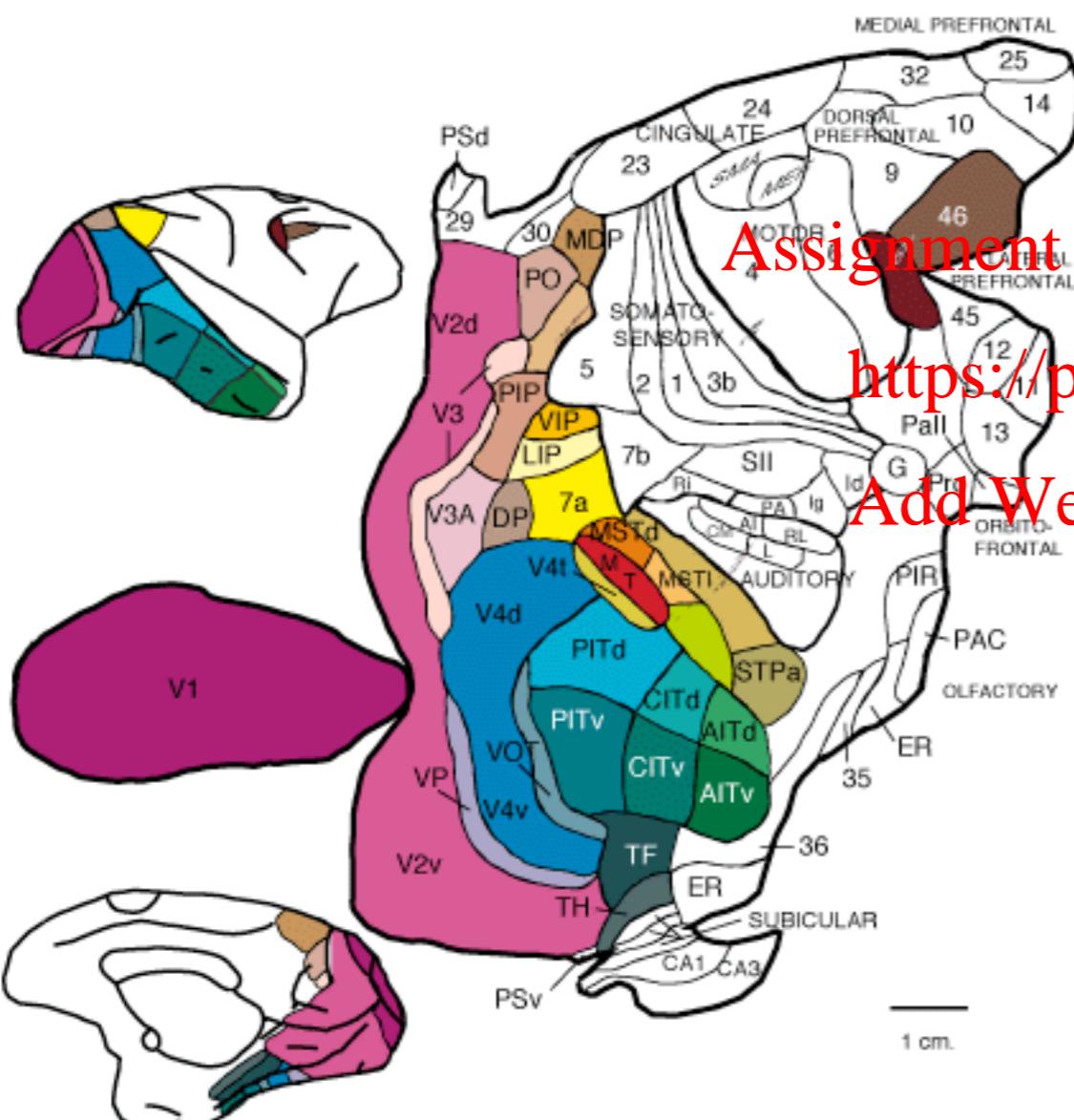
Assignment Project Exam Help

<https://powcoder.com>

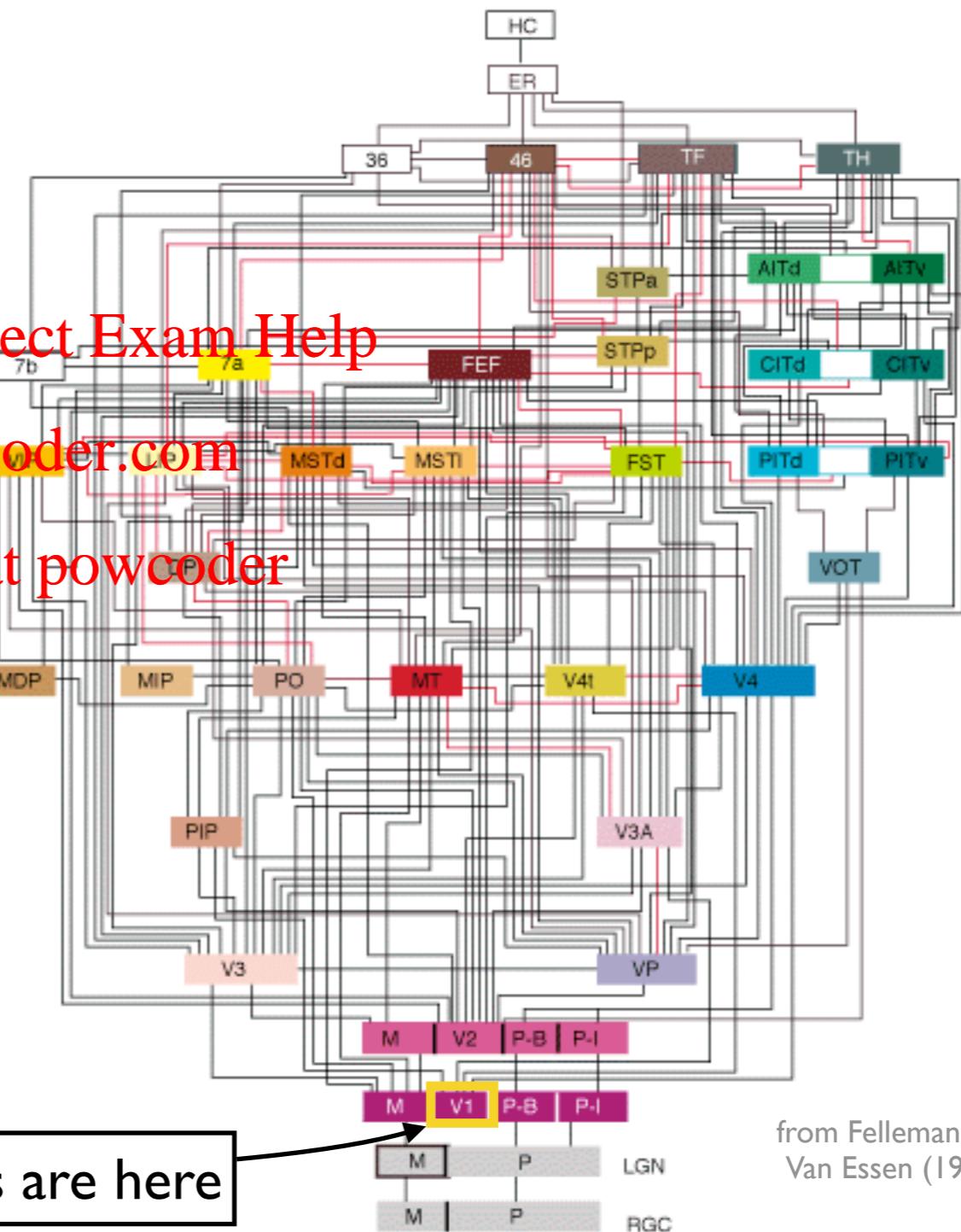
Add WeChat powcoder

Hierarchy of brain areas in the mammalian visual system

Flat map of macaque monkey brain



Hierarchy of brain areas



Assignment Project Exam Help

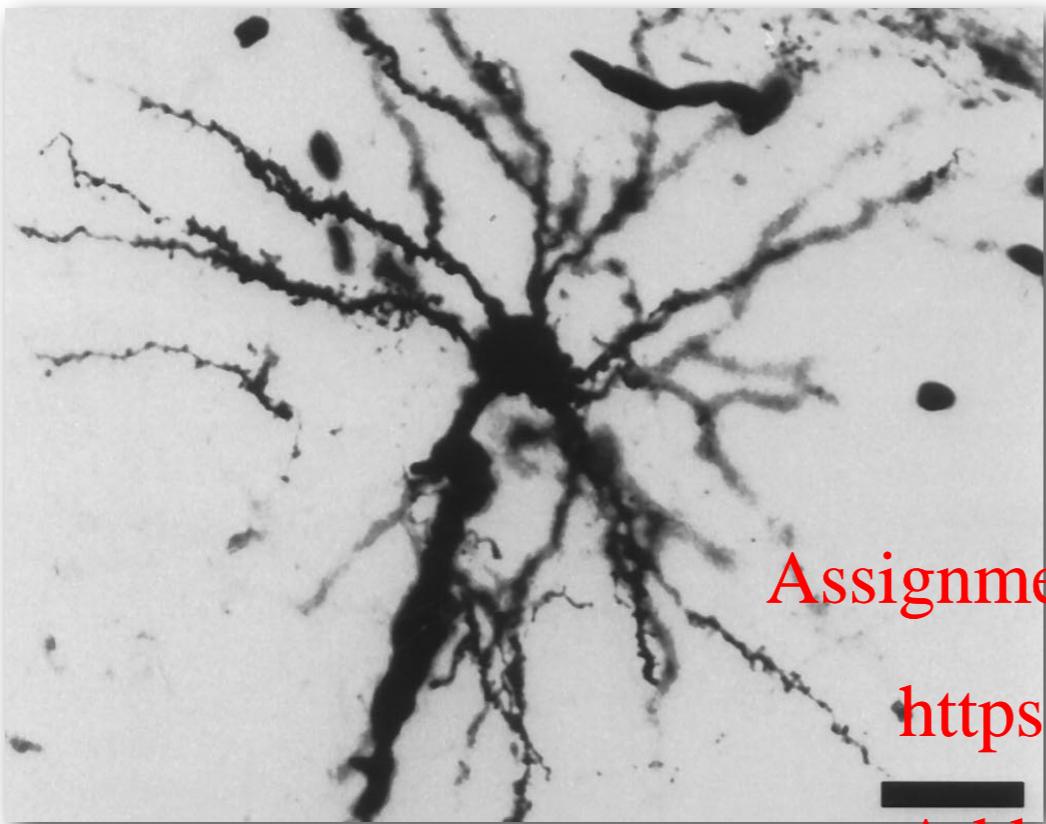
<https://powcoder.com>

Add WeChat powcoder

Simple and Complex Cells are here

from Felleman and
Van Essen (1991)

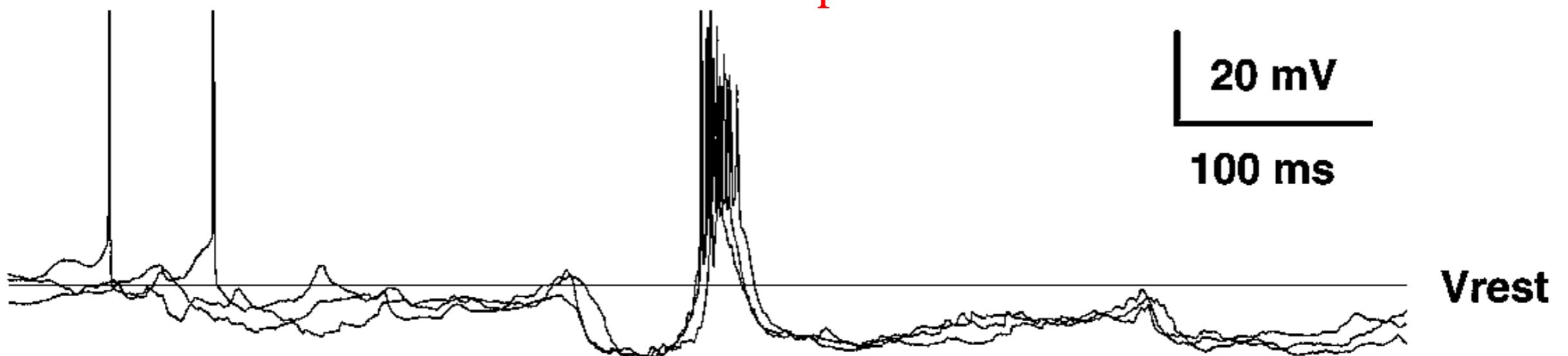
Real neurons



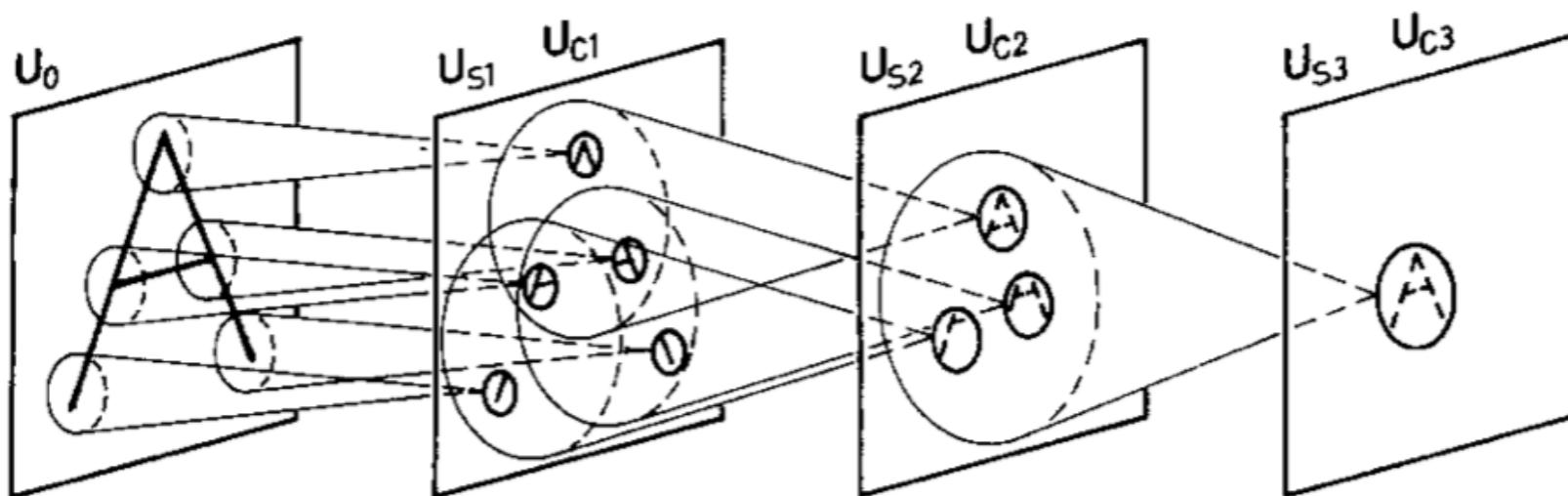
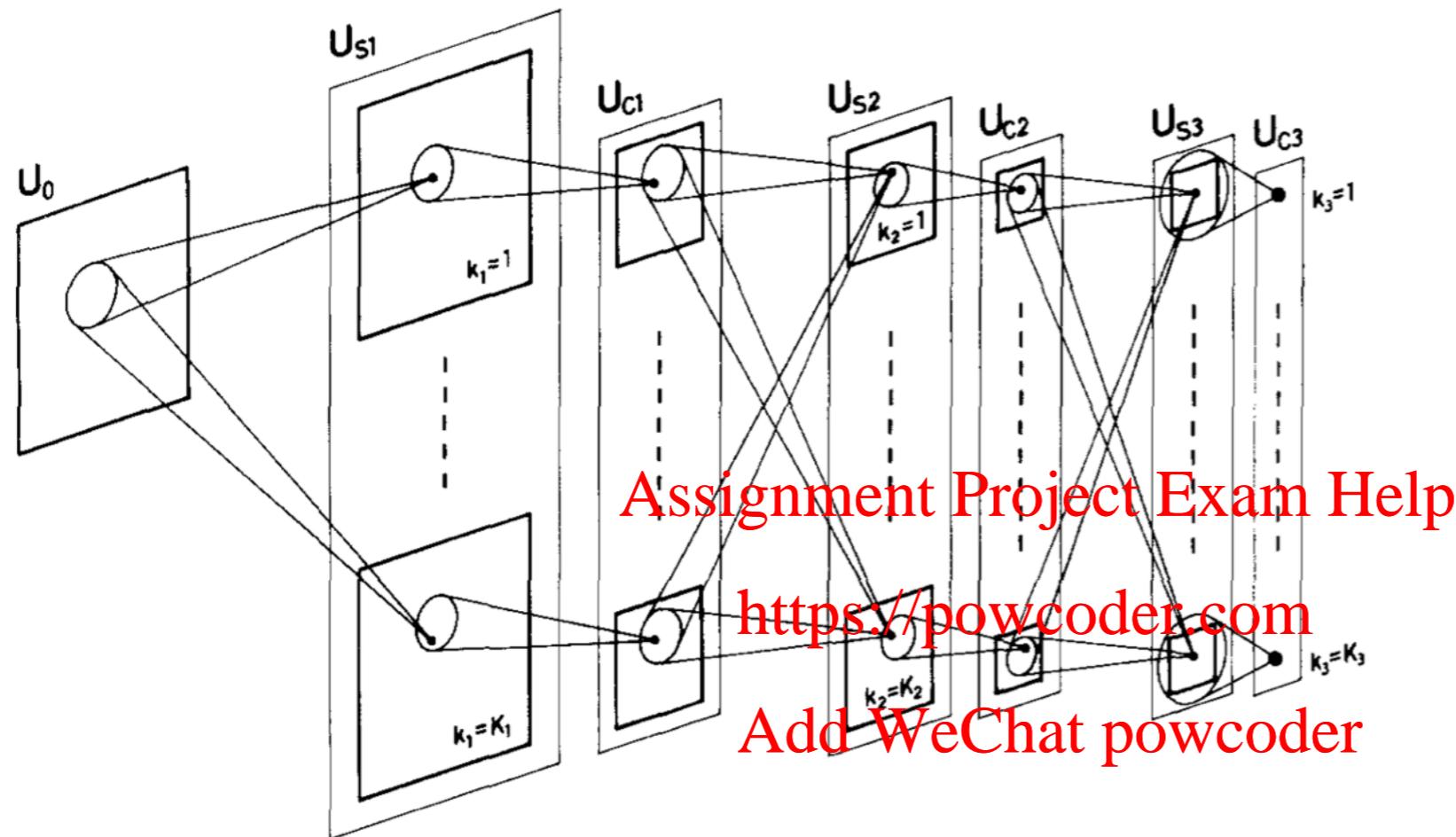
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



The Neocognitron (Fukushima, 1980)



- features are pooled from local areas
- the aim is to allow recognition of a variety of feature positions, while preserving relative position

0 0 0 1 1 1 1 1 1 2

2 2 2 2 2 2 3 3 3

Assignment Project Exam Help

3 4 4 9 9 4 5 5 5

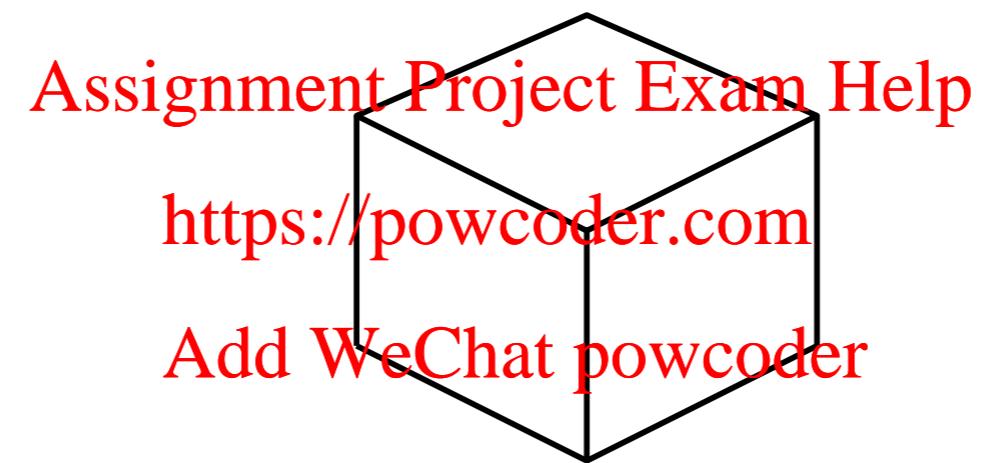
<https://powcoder.com>

Add WeChat powcoder

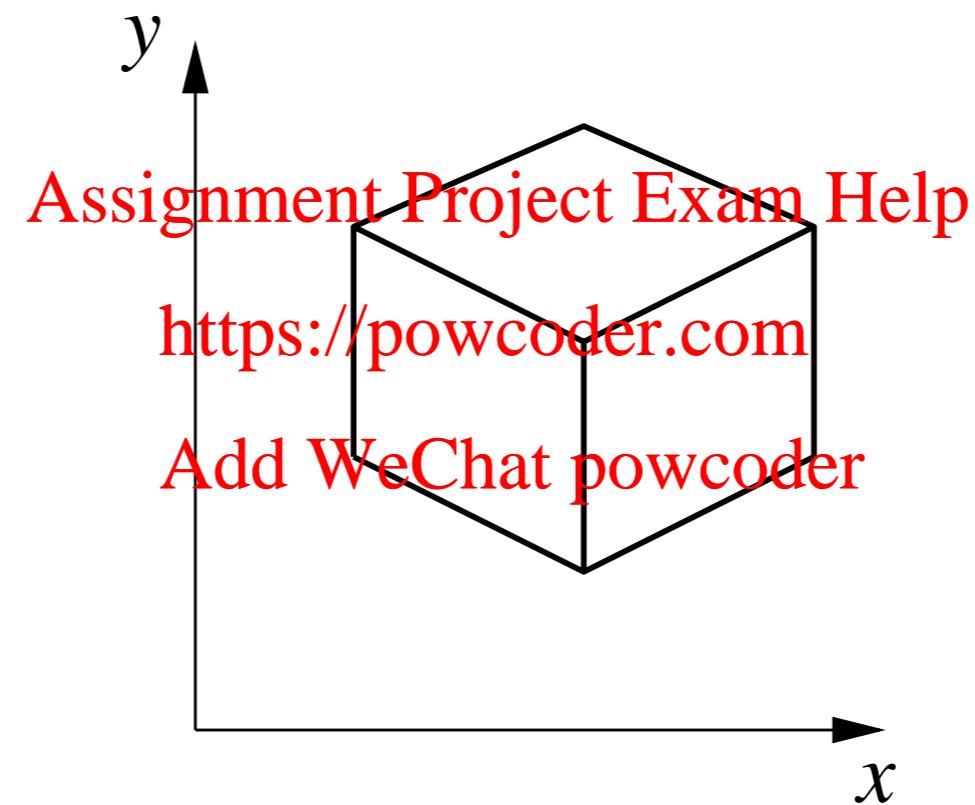
4 4 7 7 7 7 8 8 8

8 8 8 7 9 4 9 9 9

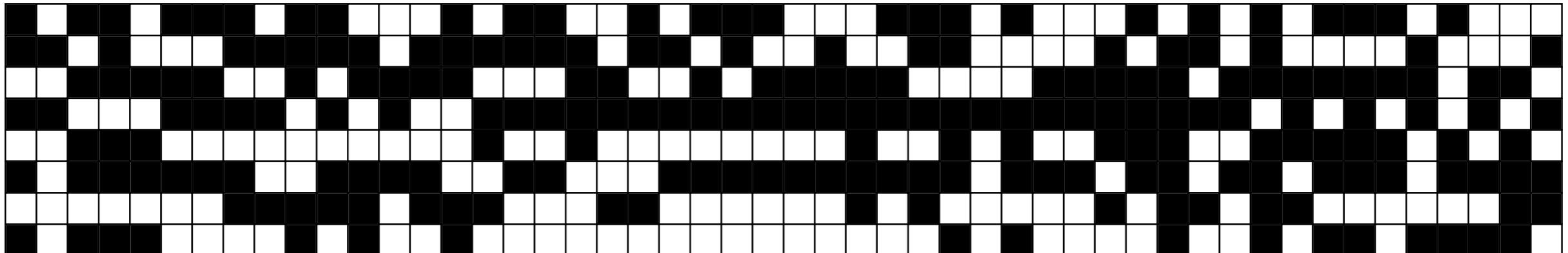
What do you see?



What do you see?



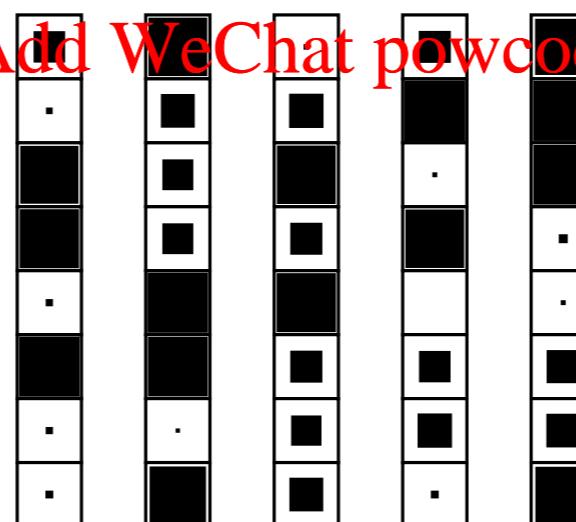
The data: a set of stochastic binary patterns



Each column is a distinct eight dimensional binary feature.

<https://powcoder.com>

Add WeChat powcoder

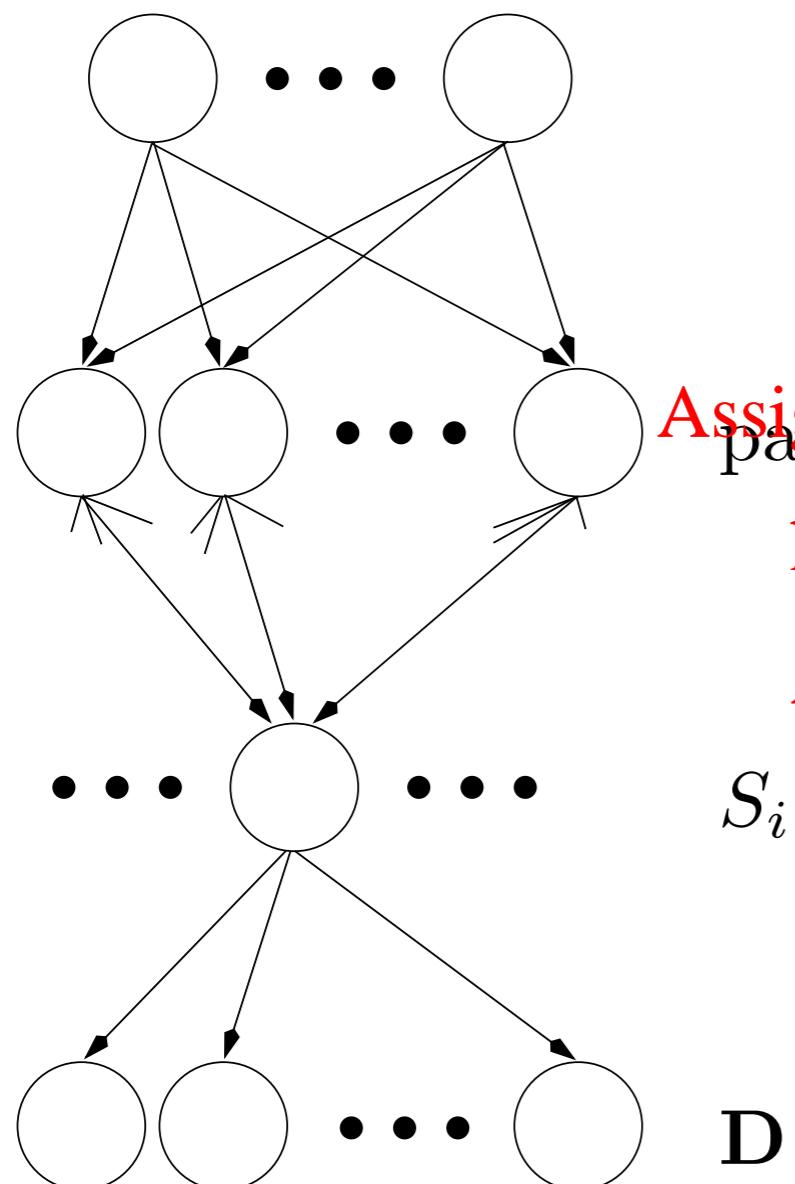


true hidden causes of the data

This is a *learning* problem, which
we'll cover in later lecture.

Hierarchical Statistical Models

A Bayesian belief network:



The joint probability of binary states is

$$P(\mathbf{S}|\mathbf{W}) = \prod_i P(S_i|\text{pa}(S_i), \mathbf{W})$$

The probability S_i depends only on its parents:

Assignment Project Exam Help
<https://powcoder.com>

$$P(S_i|\text{pa}(S_i), \mathbf{W}) = \begin{cases} h(\sum_j S_j w_{ji}) & \text{if } S_i = 1 \\ 1 - h(\sum_j S_j w_{ji}) & \text{if } S_i = 0 \end{cases}$$

Add WeChat powcoder

The function h specifies how causes are combined, $h(u) = 1 - \exp(-u)$, $u > 0$.

Main points:

- hierarchical structure allows model to form high order representations
- upper states are priors for lower states
- weights encode higher order features

Another approach: Inference by *stochastic simulation*

Basic idea:

1. Draw N samples from a sampling distribution S
2. Compute an approximate posterior probability
3. Show this converges to the true probability

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

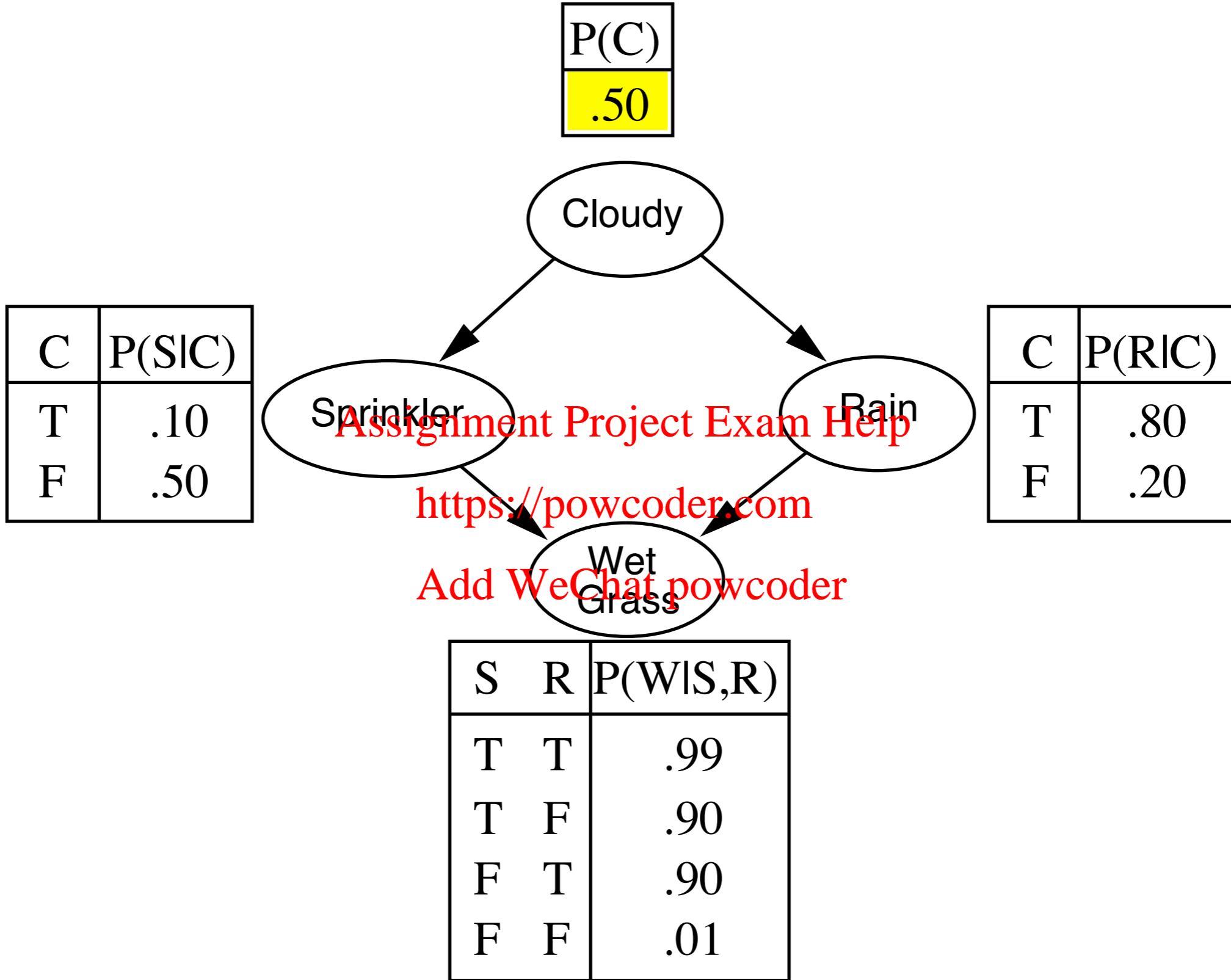
Sampling with no evidence (from the prior)

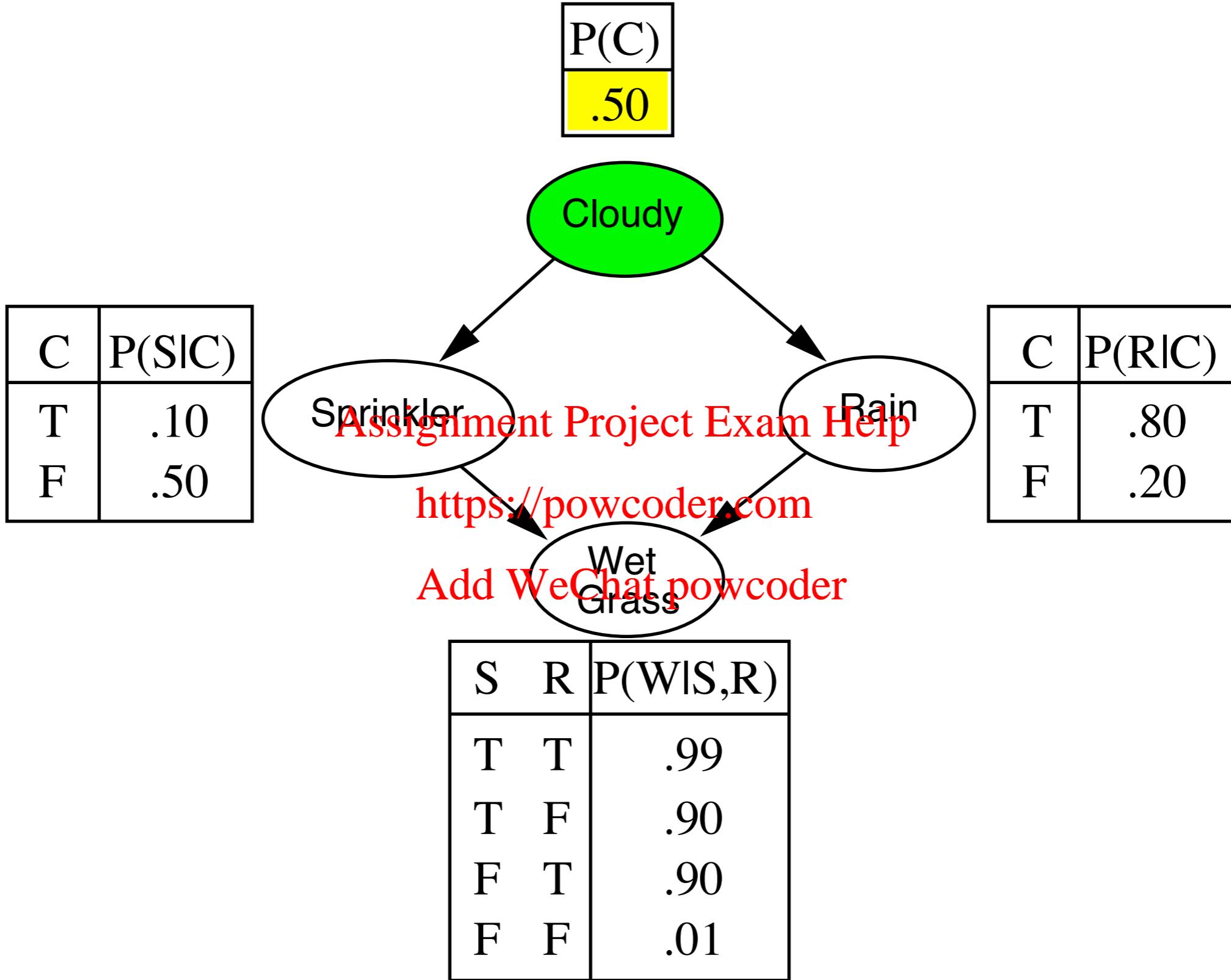
```
function PRIOR-SAMPLE(bn) returns an event sampled from bn
  inputs: bn, a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 
    x  $\leftarrow$  an event with n elements
    for i = 1 to n do
       $x_i \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
      given the values of  $\text{Parents}(X_i)$  in x
    return x
```

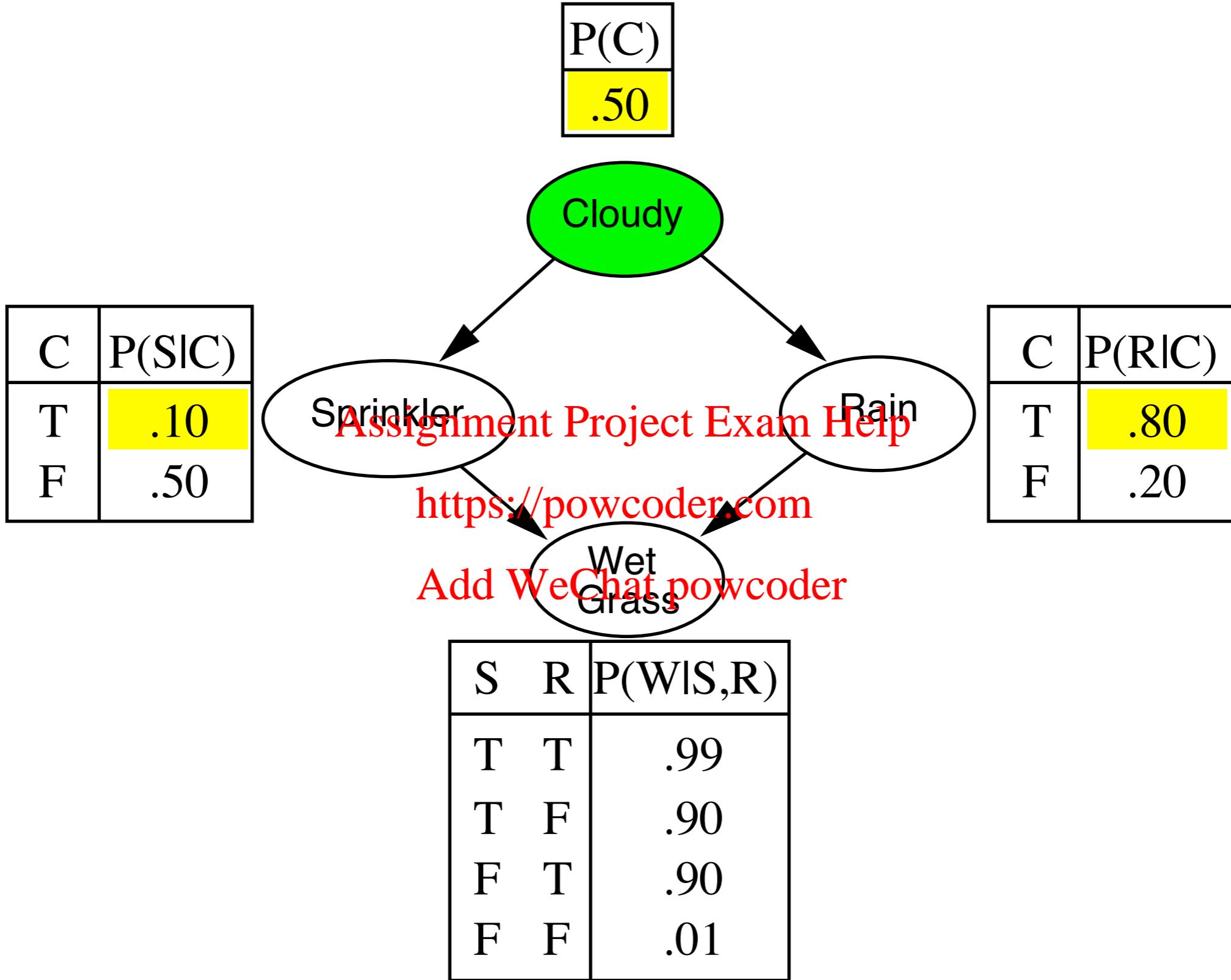
Assignment Project Exam Help

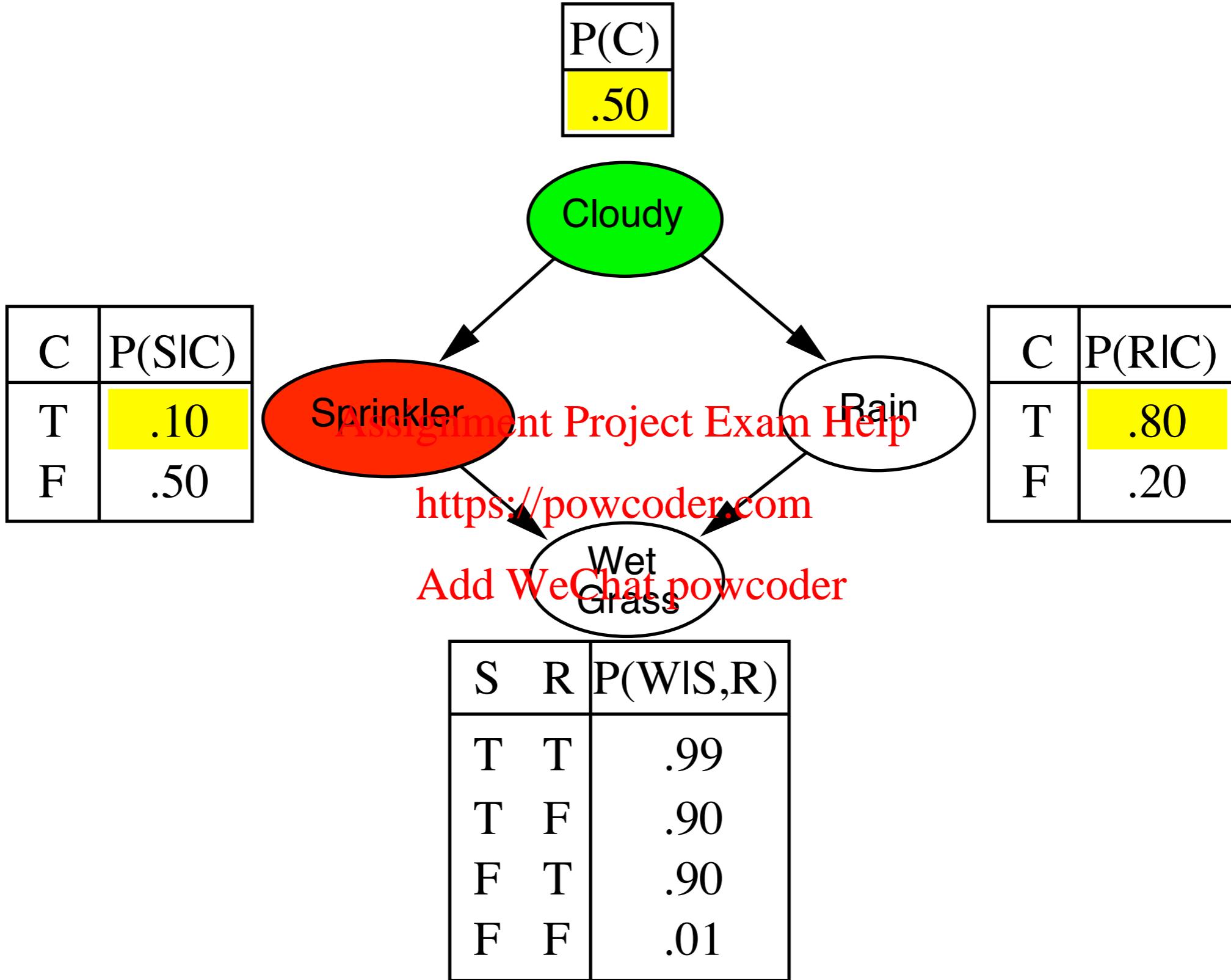
<https://powcoder.com>

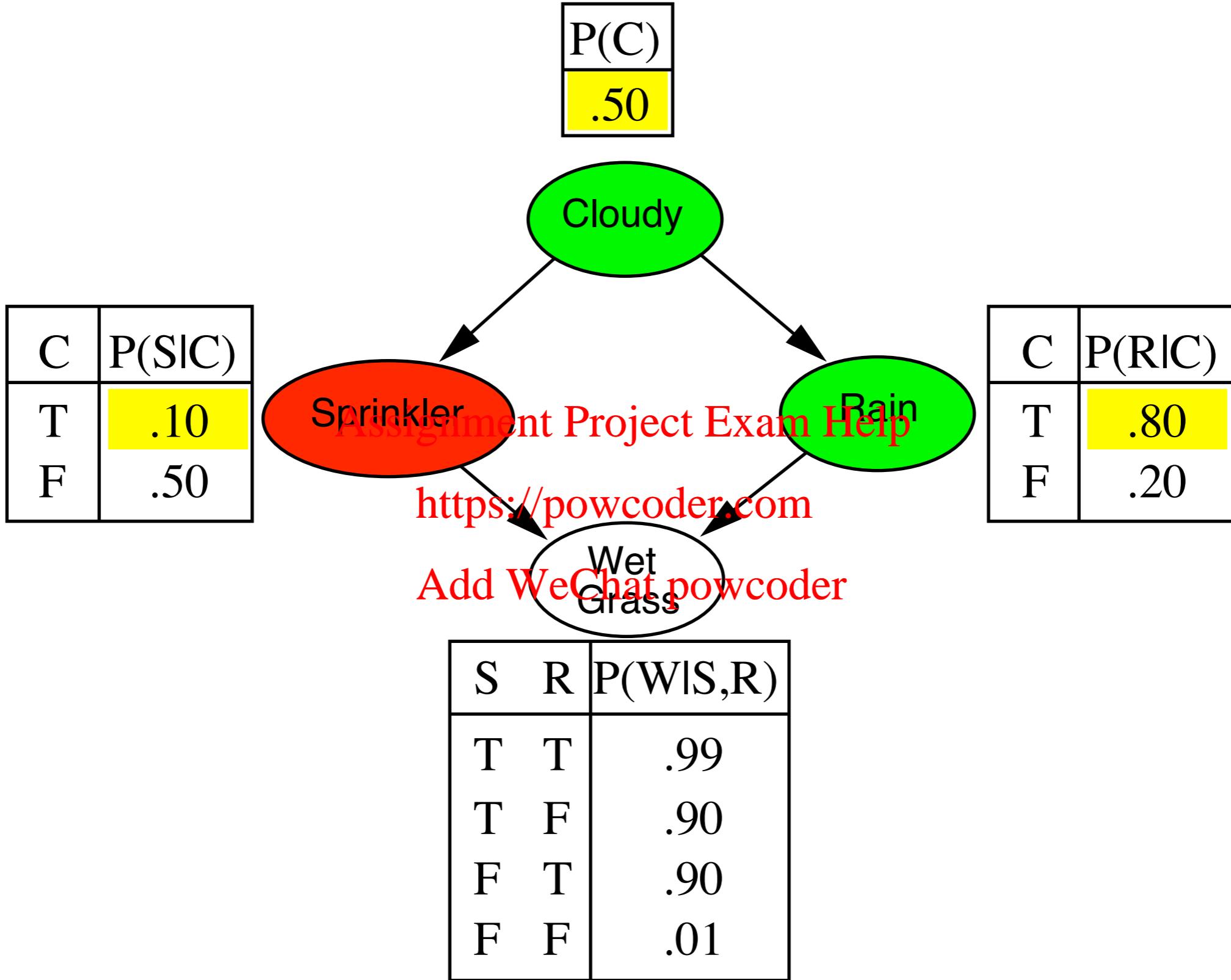
Add WeChat powcoder

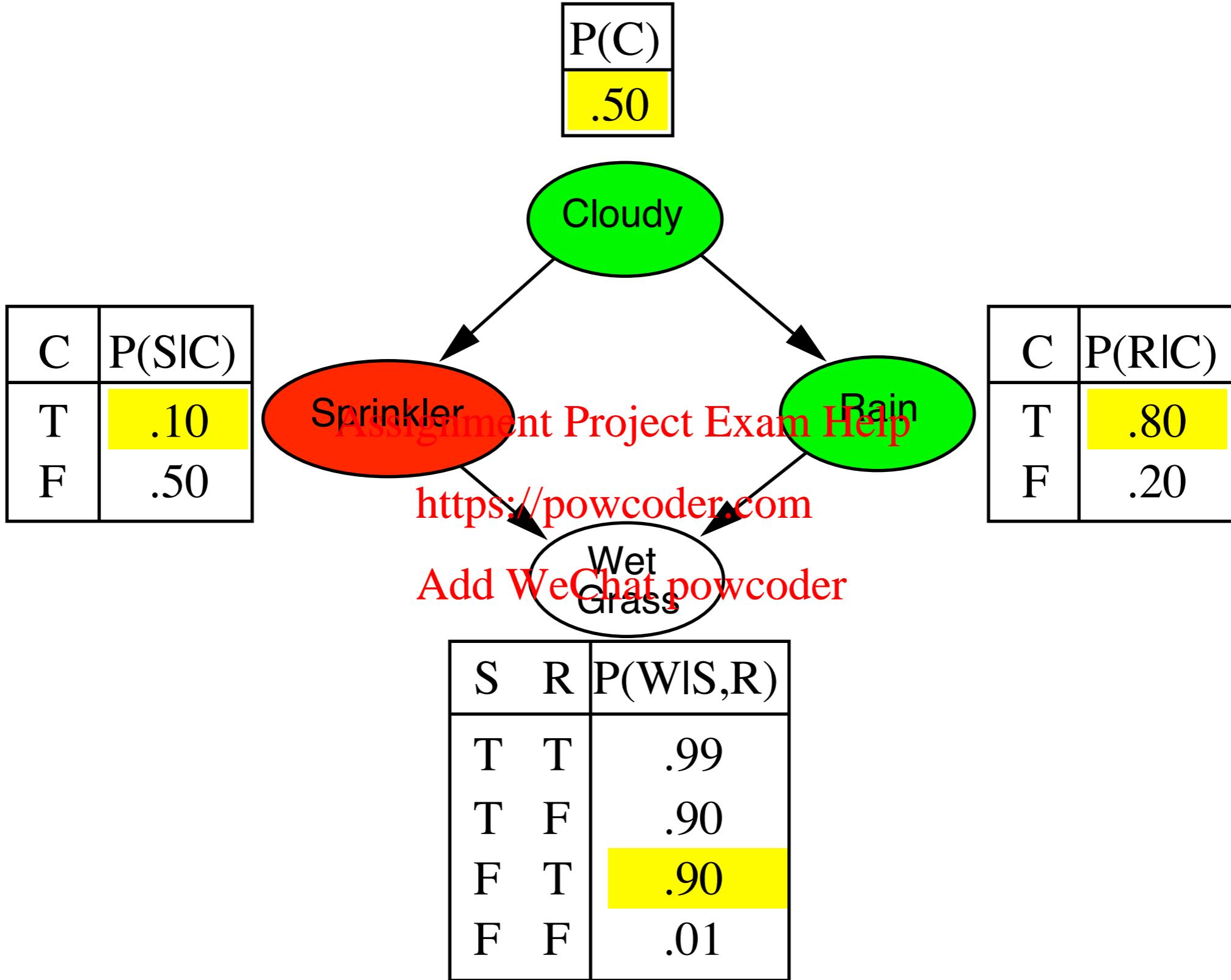


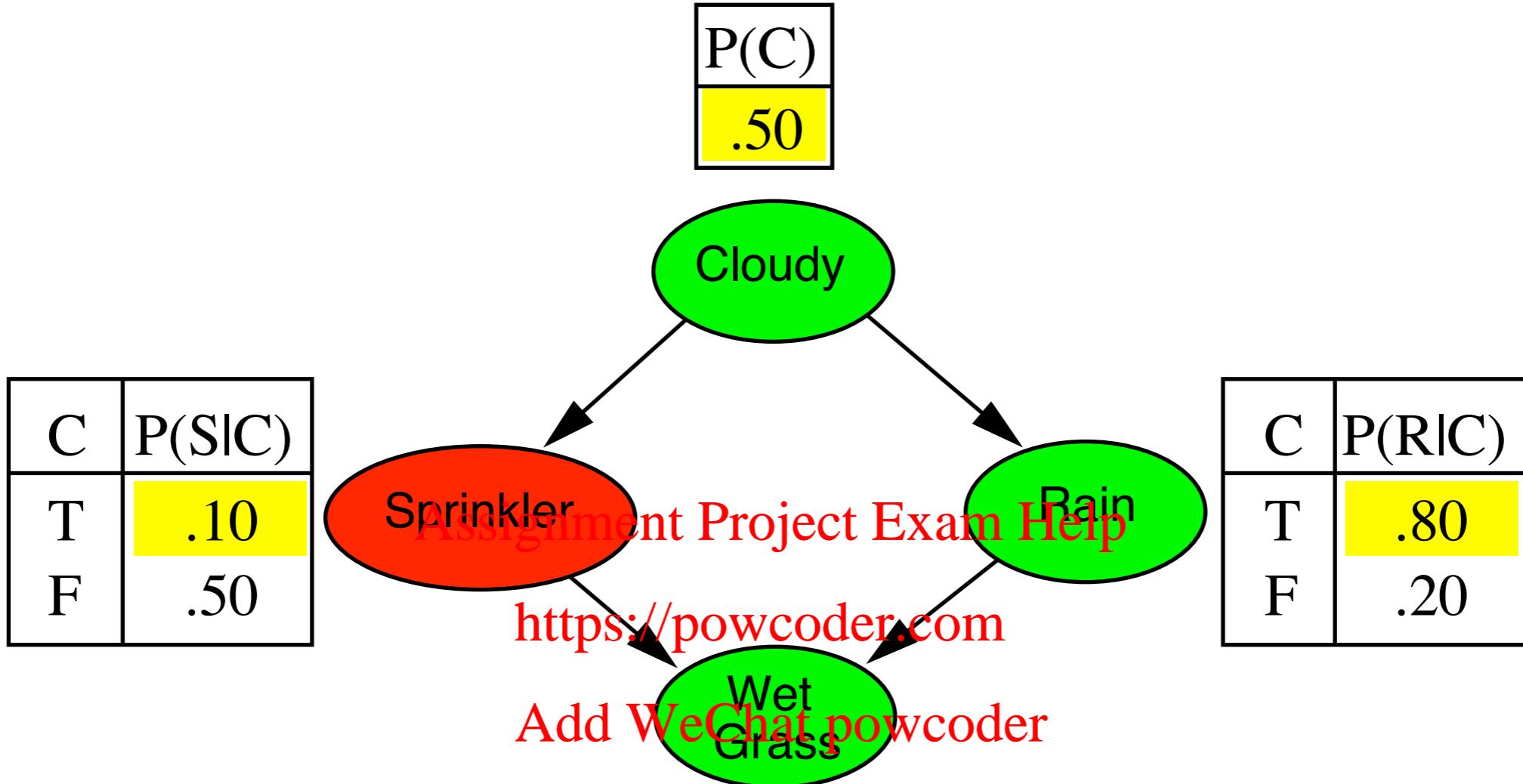












S	R	$P(W S,R)$
T	T	.99
T	F	.90
F	T	.90
F	F	.01

Keep sampling
to estimate joint
probabilities of
interest.

What if we do have some evidence? Rejection sampling.

$\hat{P}(X|e)$ estimated from samples agreeing with e

```
function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N$ , a vector of counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x \leftarrow$  PRIOR-SAMPLE( $bn$ )
    if  $x$  is consistent with  $e$  then
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $N[X]$ )  
https://powcoder.com
```

Add WeChat powcoder

E.g., estimate $P(Rain|Sprinkler = true)$ using 100 samples

27 samples have $Sprinkler = true$

Of these, 8 have $Rain = true$ and 19 have $Rain = false$.

$\hat{P}(Rain|Sprinkler = true) = \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$

Similar to a basic real-world empirical estimation procedure

Analysis of rejection sampling

$$\begin{aligned}\hat{P}(X|e) &= \alpha N_{PS}(X, e) && (\text{algorithm defn.}) \\ &= N_{PS}(X, e)/N_{PS}(e) && (\text{normalized by } N_{PS}(e)) \\ &\approx P(X, e)/P(e) && (\text{property of PRIORSAMPLE}) \\ &= P(X|e) && (\text{defn. of conditional probability})\end{aligned}$$

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(e)$ is small

$P(e)$ drops off exponentially with number of evidence variables!

<https://powcoder.com>
Add WeChat powcoder

Approximate inference using Markov Chain Monte Carlo (MCMC)

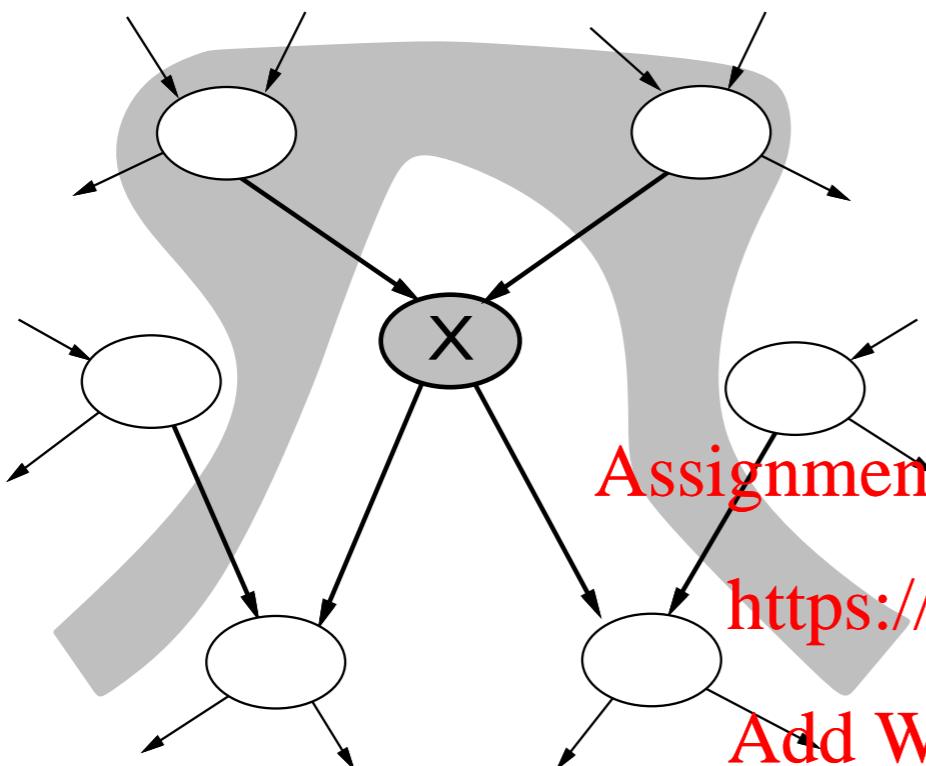
“State” of network = current assignment to all variables.

Generate next state by sampling one variable given Markov blanket
Sample each variable in turn, keeping evidence fixed

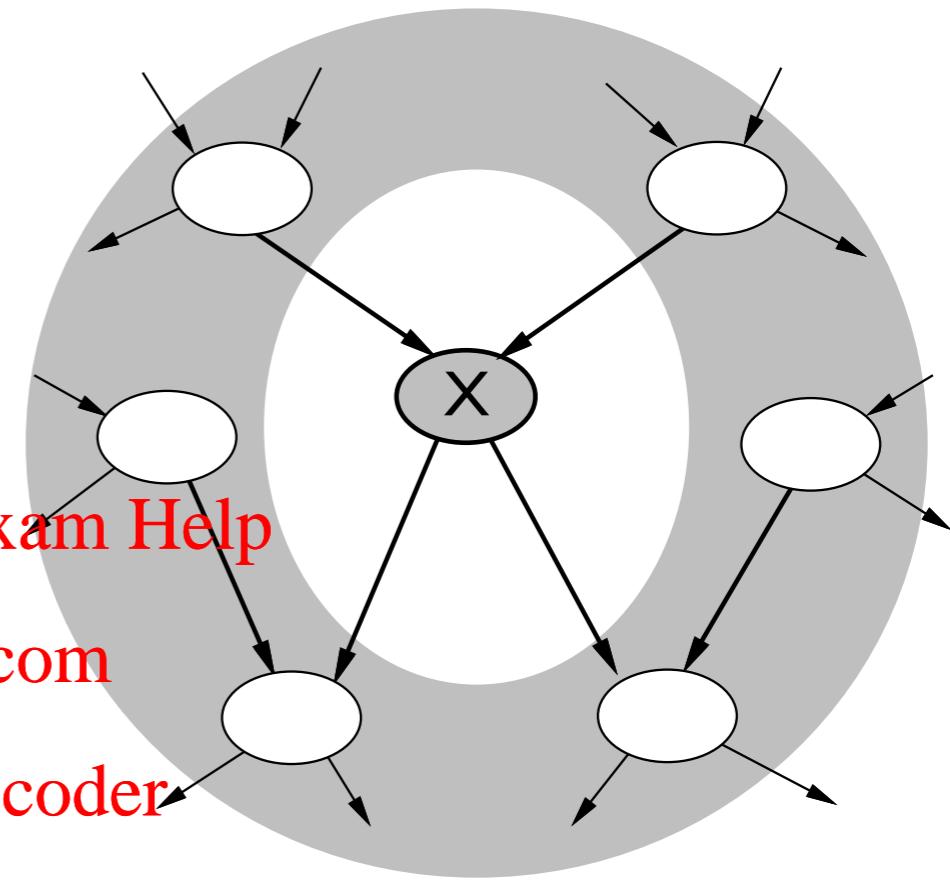
```
function MCMC-ASK( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $P(X|\mathbf{e})$ 
    local variables:  $\mathbf{N}[X]$ , a vector of counts over  $X$ , initially zero
         $Z$ , The evidence variables in  $bn$ 
         $\mathbf{x}$ , the current state of the network, initially copied from  $\mathbf{e}$ 
        https://powcoder.com
    initialize  $\mathbf{x}$  with random values for the variables in  $Y$ 
    for  $j = 1$  to  $N$  do          Add WeChat powcoder
        for each  $Z_i$  in  $Z$  do
            sample the value of  $Z_i$  in  $\mathbf{x}$  from  $P(Z_i|mb(Z_i))$ 
            given the values of  $MB(Z_i)$  in  $\mathbf{x}$ 
             $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
    return NORMALIZE( $\mathbf{N}[X]$ )
```

Can also choose a variable to sample at random each time

The extent of dependencies in Bayesian networks



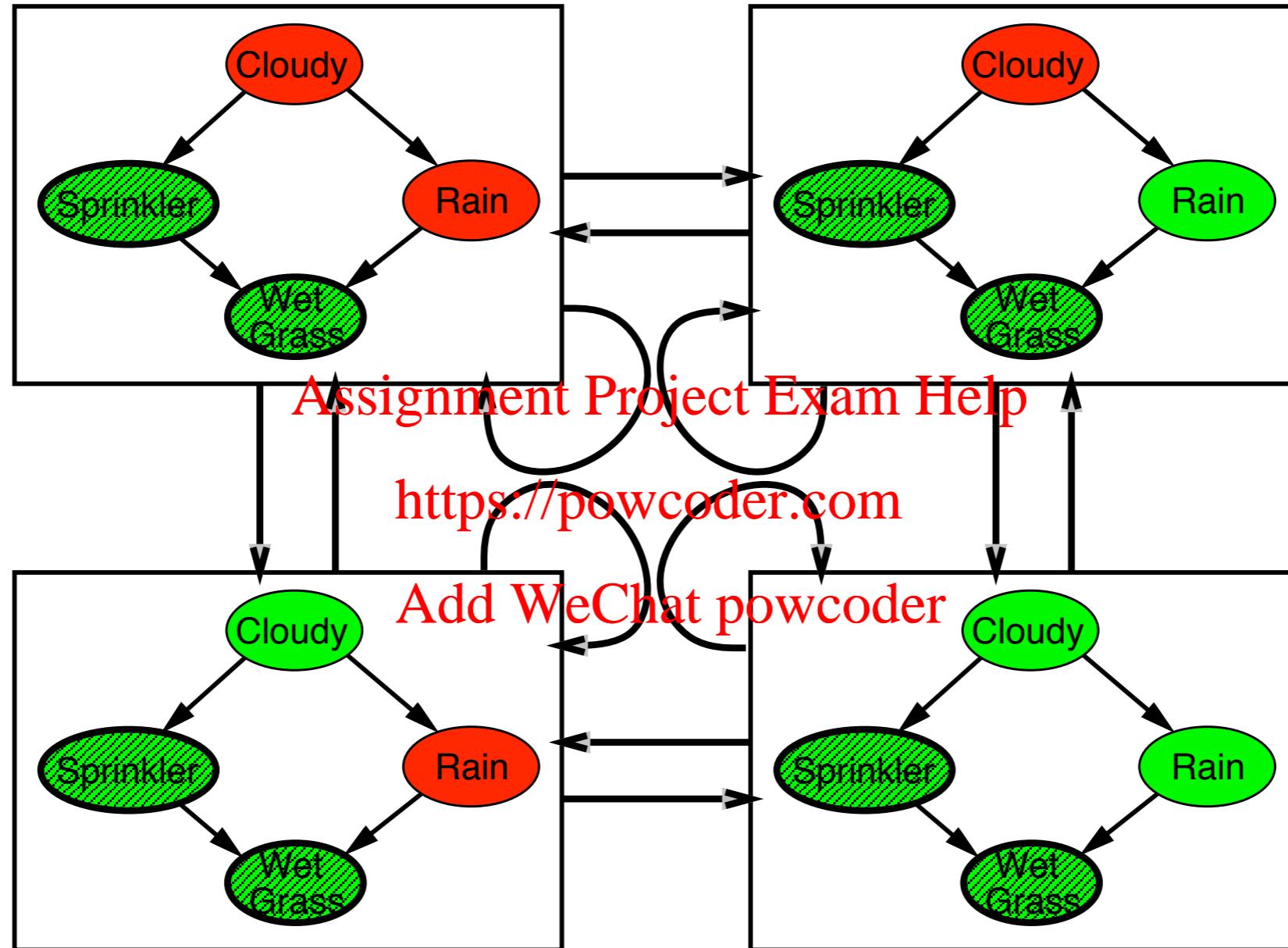
Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder



- A node X is conditionally independent of its non-descendants given its parents
- A node X is conditionally independent of all the other nodes in the network given its **Markov blanket**.

The Markov chain

With $\text{Sprinkler} = \text{true}$, $\text{WetGrass} = \text{true}$, there are four states:



Wander about for a while, average what you see

After obtaining the MCMC samples

Estimate $P(Rain | Sprinkler = \text{true}, WetGrass = \text{true})$

Sample *Cloudy* or *Rain* given its Markov blanket, repeat.

Count number of times *Rain* is true and false in the samples.

E.g., visit 100 states

31 have $\text{Rain} = \text{true}$, 69 have $\text{Rain} = \text{false}$

$$\hat{P}(Rain|Sprinkler=true, WetGrass=true) = \text{NORMALIZE}(\langle 31, 69 \rangle) \equiv \langle 0.31, 0.69 \rangle$$

Theorem: chain approaches stationary distribution: [Add weChat powcoder](#)

long-run fraction of time spent in each state is exactly proportional to its posterior probability

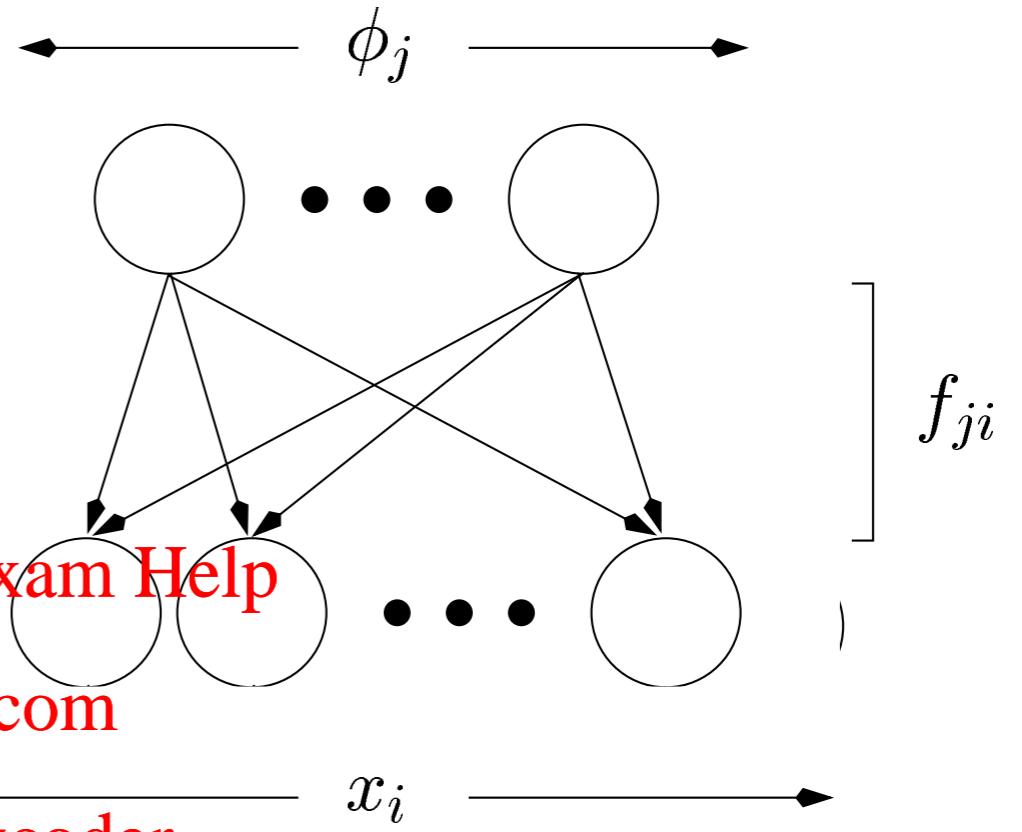
But:

- I. Difficult to tell when samples have converged. Theorem only applies in limit, and it could take time to “settle in”.
 - 2. Can also be inefficient if each state depends on many other variables.

Gibbs sampling (back to the noisy-OR example)

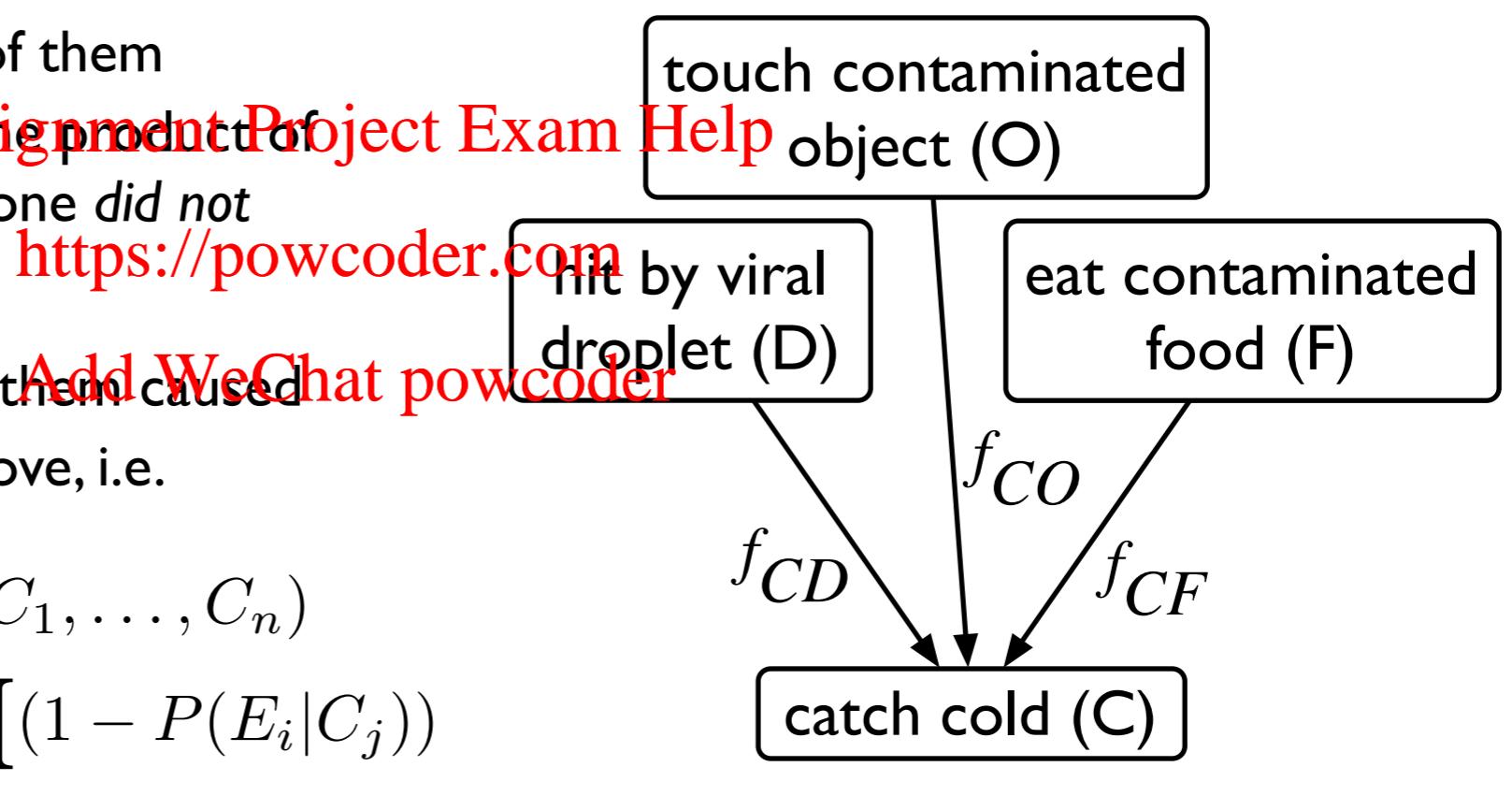
- Model represents stochastic binary features.
- Each input x_i encodes the probability that the i th binary input feature is present.
- The set of features represented by ϕ_j is defined by weights f_{ji} , which encode the probability that feature i is an instance of ϕ_j .
- Trick: It's easier to adapt weights in an unbounded space, so use the transformation:

$$f = 1/(1 + \exp(-w))$$



Beyond tables: modeling causal relationships using Noisy-OR

- We assume each cause C_j can produce effect E_i with probability f_{ij} .
- The noisy-OR model assumes the parent causes of effect E_i contribute independently.
- The probability that none of them caused effect E_i is simply the product of the probabilities that each one did not cause E_i .

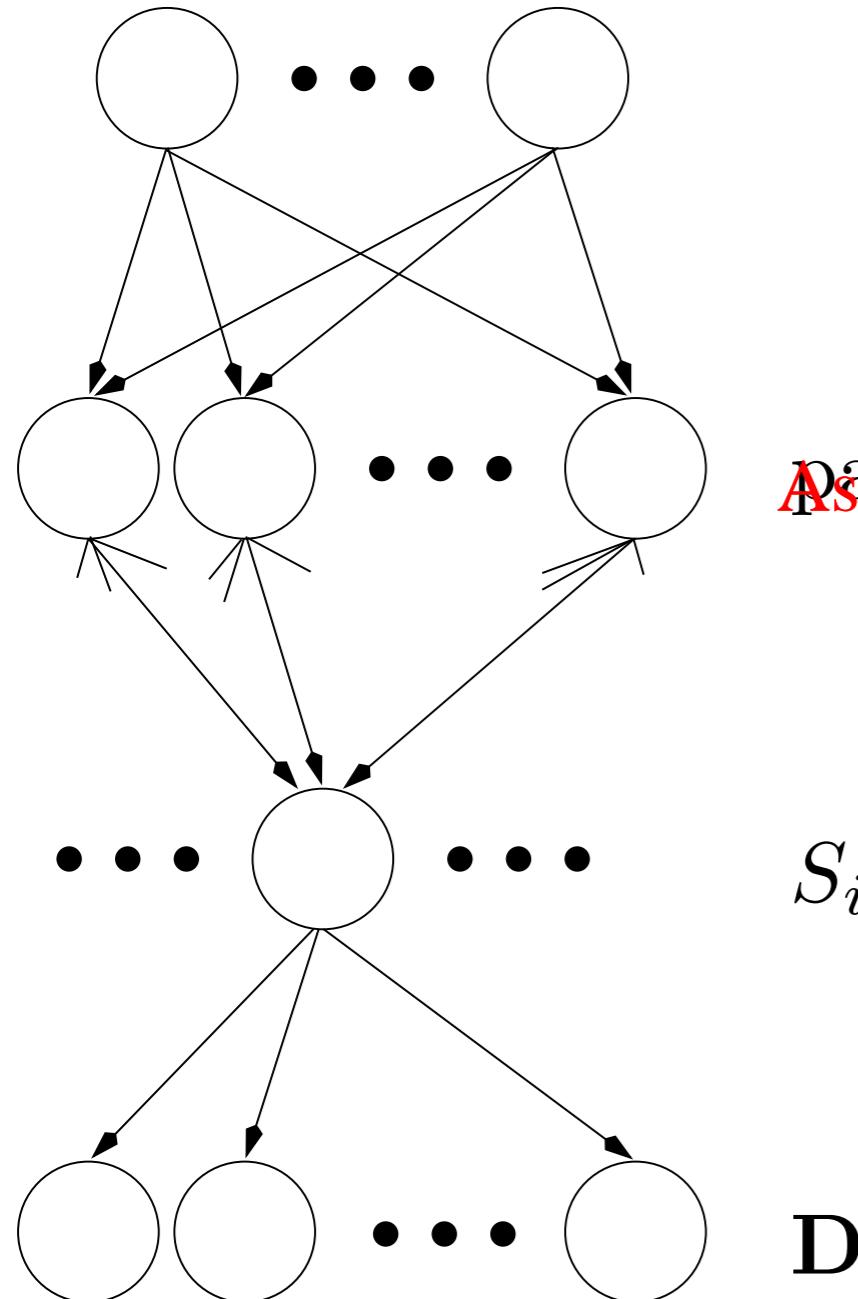


$$\begin{aligned}
 P(E_i | \text{par}(E_i)) &= P(E_i | C_1, \dots, C_n) \\
 &= 1 - \prod_i (1 - P(E_i | C_j)) \\
 &= 1 - \prod_i (1 - f_{ij})
 \end{aligned}$$

$$\begin{aligned}
 P(C|D, O, F) &= \\
 &1 - (1 - f_{CD})(1 - f_{CO})(1 - f_{CF})
 \end{aligned}$$

Hierarchical Statistical Models

A Bayesian belief network:



The joint probability of binary states is

$$P(\mathbf{S}|\mathbf{W}) = \prod_i P(S_i|\text{pa}(S_i), \mathbf{W})$$

The probability S_i depends only on its parents:

$$P(S_i|\text{pa}(S_i), \mathbf{W}) =$$

Assignment Project Exam Help
<https://powcoder.com>

$$\begin{cases} h(\sum_j S_j w_{ji}) & \text{if } S_i = 1 \\ 1 - h(\sum_j S_j w_{ji}) & \text{if } S_i = 0 \end{cases}$$

Add WeChat powcoder

The function h specifies how causes are combined, $h(u) = 1 - \exp(-u)$, $u > 0$.

Main points:

- hierarchical structure allows model to form high order representations
- upper states are priors for lower states
- weights encode higher order features

Inferring the best representation of the observed variables

- Given on the input \mathbf{D} , there is no simple way to determine which states are the input's most likely causes.
 - Computing the most probable network state is an *inference* process
 - we want to find the explanation of the data with highest probability
 - this can be done efficiently with *Gibbs sampling*
- Assignment Project Exam Help**
- Gibbs sampling is another example of an MCMC method
- Key idea:

<https://powcoder.com>

The samples are guaranteed to converge to the true posterior probability distribution

Gibbs Sampling

Gibbs sampling is a way to select an ensemble of states that are representative of the posterior distribution $P(\mathbf{S}|\mathbf{D}, \mathbf{W})$.

- Each state of the network is updated iteratively according to the probability of S_i given the remaining states.
- this conditional probability can be computed using (Neal, 1992)

Assignment Project Exam Help

$$P(S_i = a | S_j : j \neq i, \mathbf{W}) \propto P(S_i = a | \text{pa}(S_i), \mathbf{W}) \prod_{j \in \text{ch}(S_i)} P(S_j | \text{pa}(S_j), S_i = a, \mathbf{W})$$

<https://powcoder.com>
Add WeChat powcoder

- limiting ensemble of states will be typical samples from $P(\mathbf{S}|\mathbf{D}, \mathbf{W})$
- also works if any subset of states are fixed and the rest are sampled

The Gibbs sampling equations (derivation omitted)

The probability of S_i changing state given the remaining states is

$$P(S_i = 1 - S_i | S_j : j \neq i, \mathbf{W}) = \frac{1}{1 + \exp(-\Delta x_i)}$$

Δx_i indicates how much changing the state S_i changes the probability of the whole network state

Assignment Project Exam Help

<https://powcoder.com>

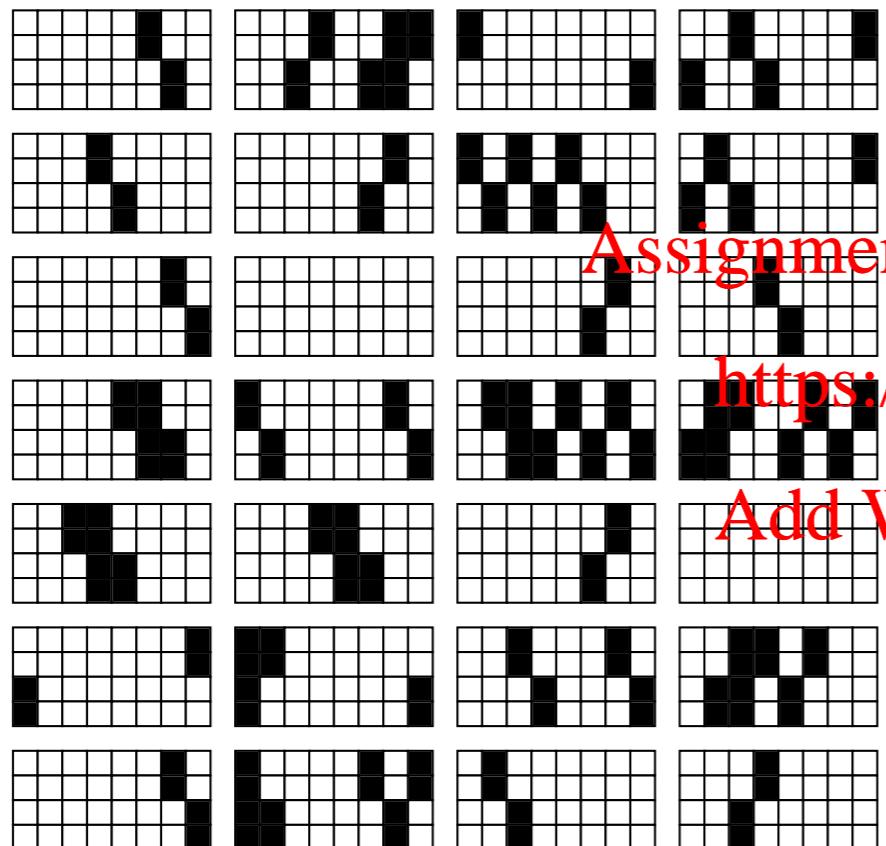
$$\begin{aligned} \Delta x_i &= \log h(u_i, 1 - S_i) - \log h(u_i, S_i) \\ &\quad + \sum_{j \in \text{ch}(S_i)} \log h(u_j + \delta_{ij}; S_j) - \log h(u_j; S_j) \end{aligned}$$

- u_i is the causal input to S_i , $u_i = \sum_k S_k w_{ki}$
- δ_j specifies the change in u_j for a change in S_i ,
 $\delta_{ij} = +S_j w_{ij}$ if $S_i = 0$, or $-S_j w_{ij}$ if $S_i = 1$

Interpretation of the Gibbs sampling equation

- The Gibbs equation can be interpreted as: $feedback + \sum feedforward$
- $feedback$: how consistent is S_i with current causes?
- $\sum feedforward$: how likely is S_i a cause of its children
- feedback allows the lower-level units to use information only
[Assignment Project Exam Help](#)
computable at higher levels <https://powcoder.com>
- feedback determines (disambiguates) the state when the
[Add WeChat powcoder](#)
feedforward input is ambiguous

The Shifter Problem

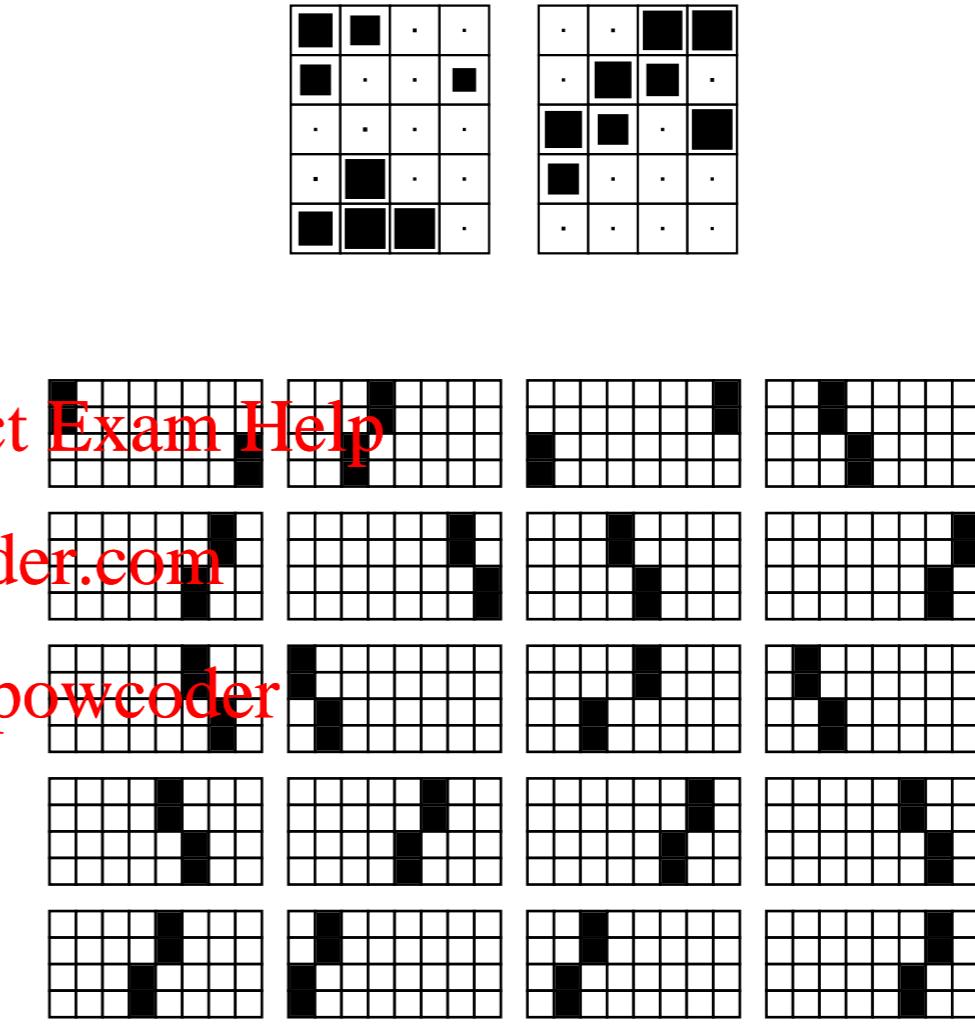


Shift patterns

Assignment Project Exam Help

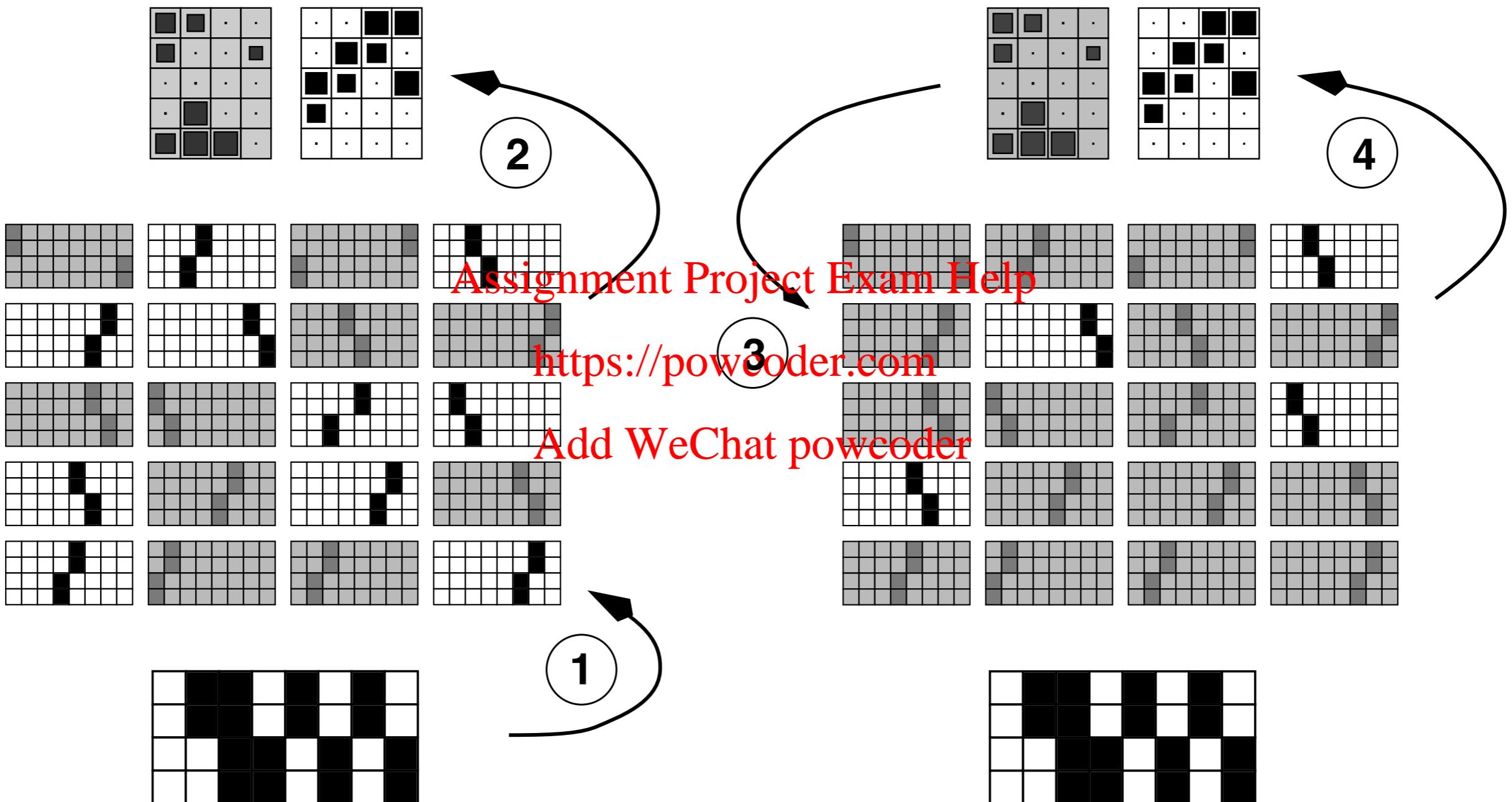
<https://powcoder.com>

Add WeChat powcoder



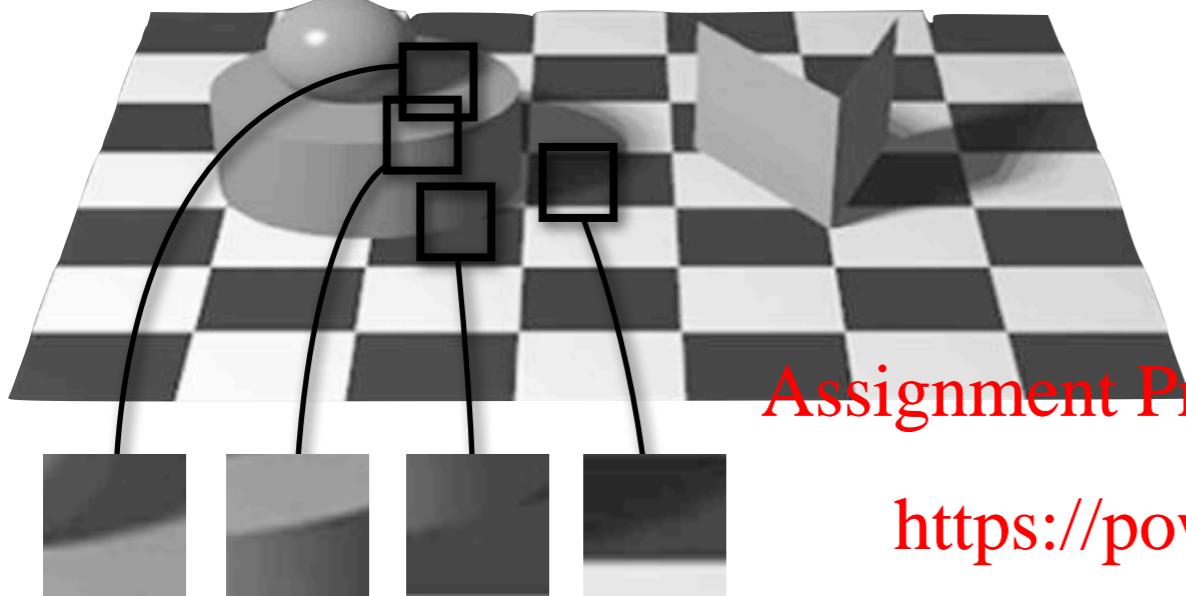
weights of a 32-20-2 network after learning

Gibbs sampling: feedback disambiguates lower-level states



Once the structure learned, the Gibbs updating converges in two sweeps.

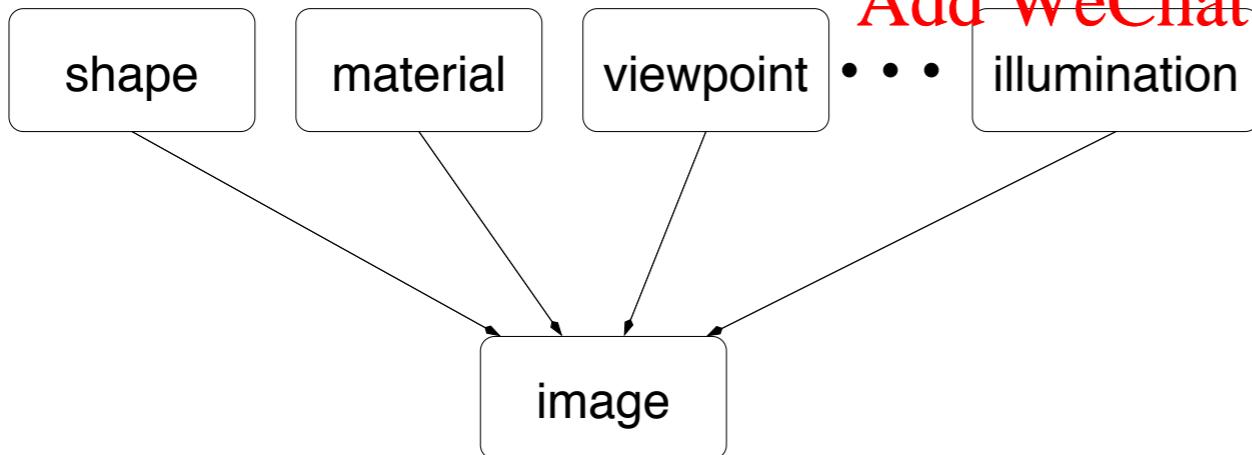
Difficulties with identifying the causal structure



Assignment Project Exam Help

<https://powcoder.com>

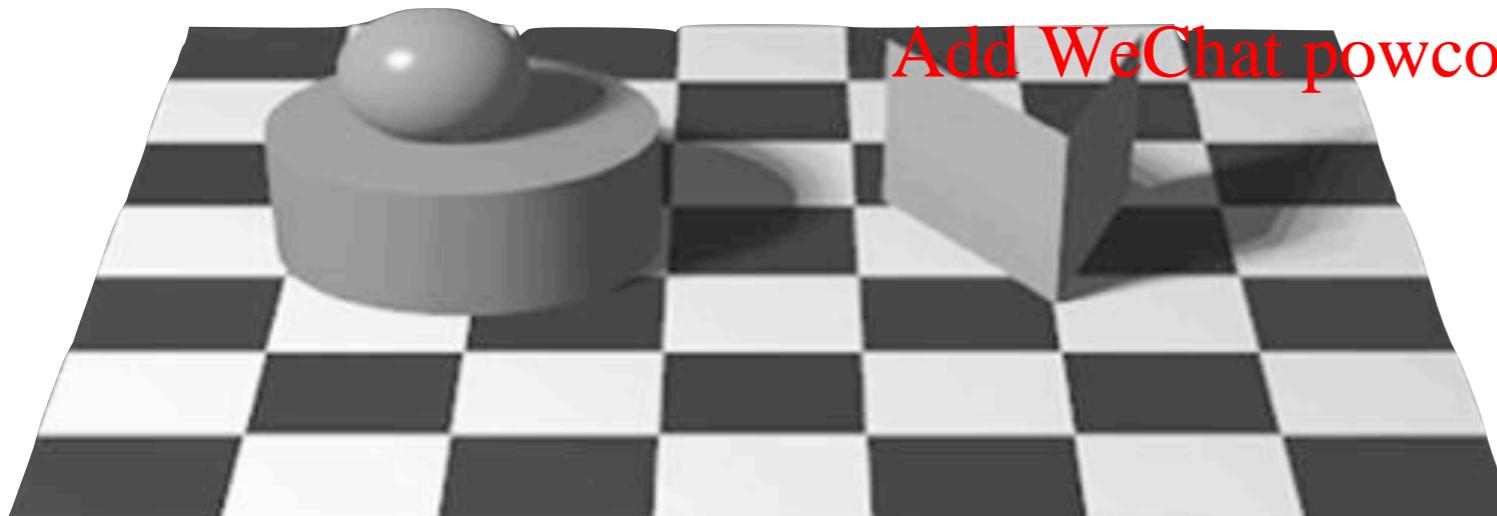
Add WeChat powcoder



- Space of S, C is huge
- If we define $P(I|S,C)$, must be able to invert it or search it
- Need efficient algorithms
 - Many unknowns: identity of objects, types of scene elements, illuminations
- Might never have encountered some structures
- Is it even the right approach?
- Can we solve a simple case?

Motivation: learn hierarchical, context dependent representations

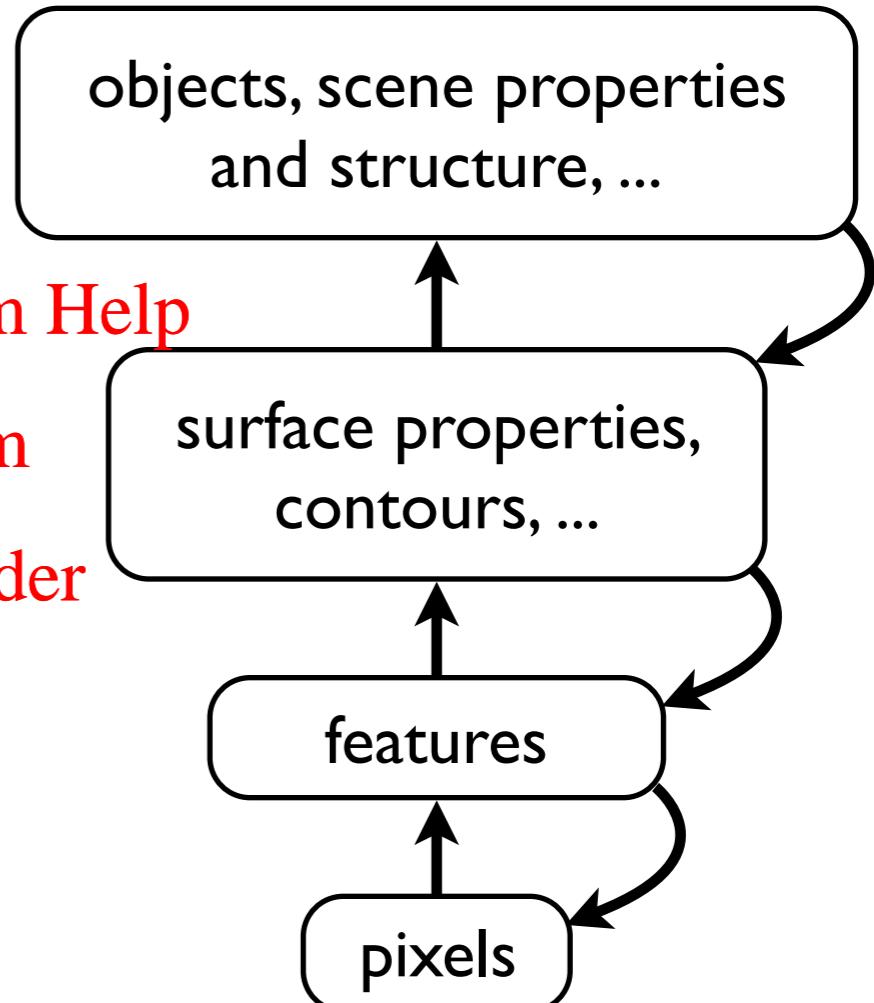
- many real-world patterns are hierarchical in structure
- interpretation of patterns depends on context
- essential for complex recognition tasks



Assignment Project Exam Help

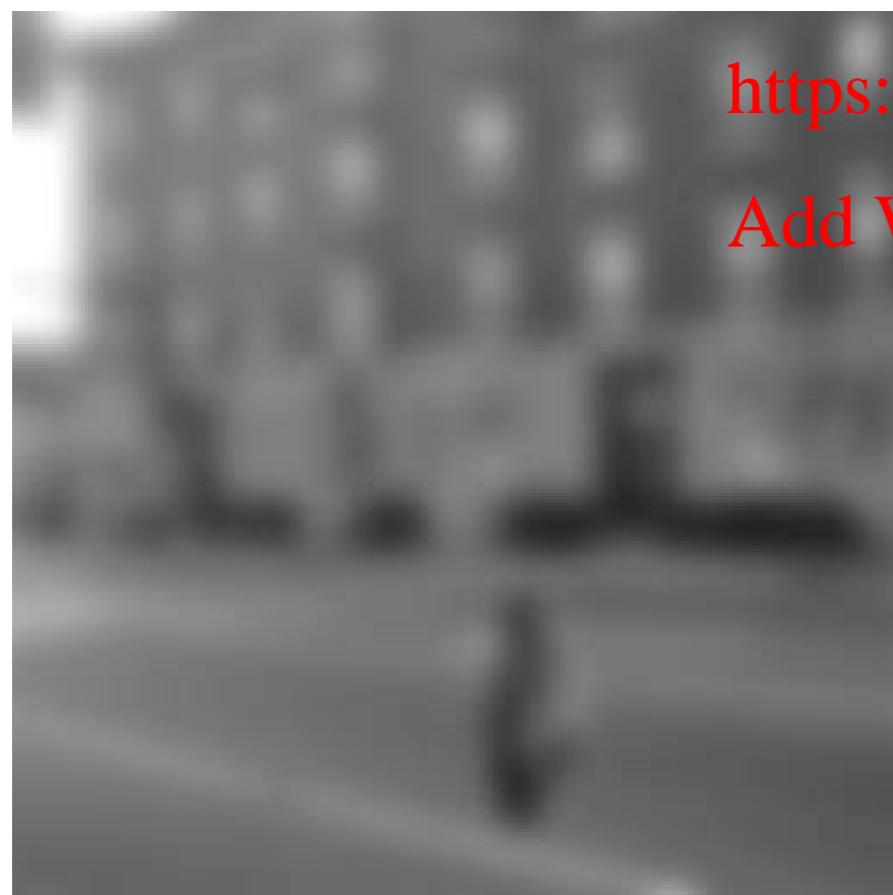
<https://powcoder.com>

Add WeChat powcoder



Motivation: hierarchical, context dependent representations

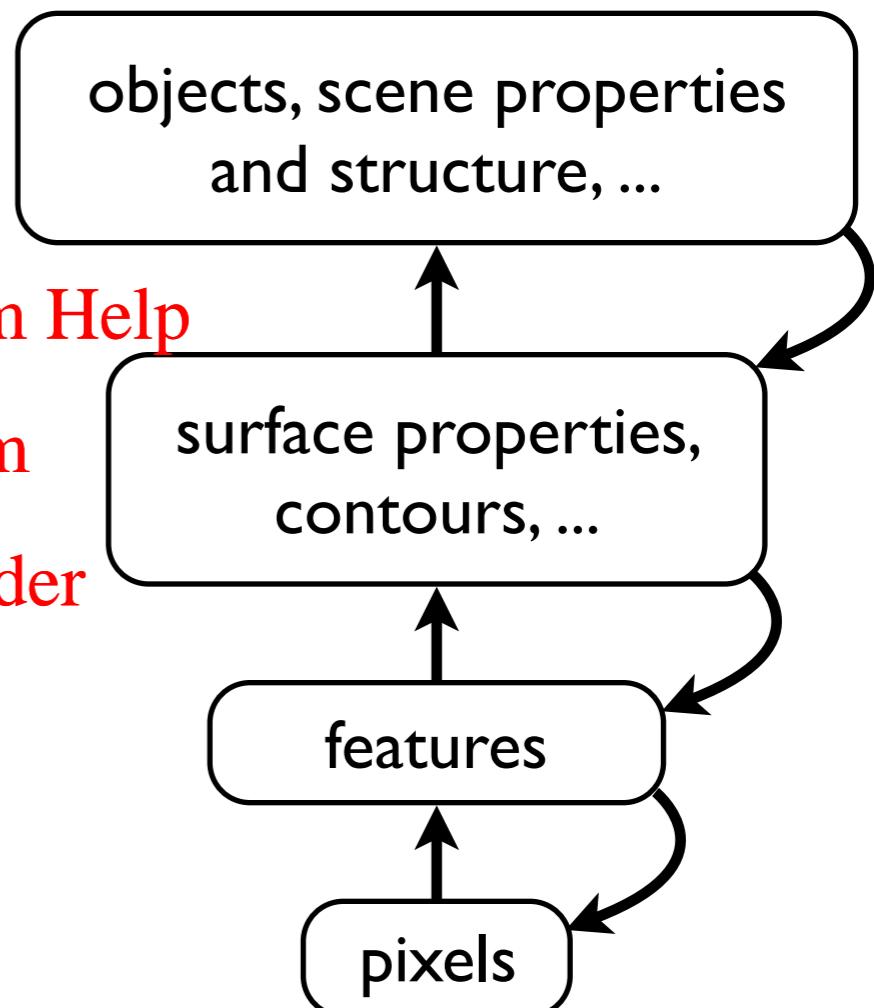
- many real-world patterns are hierarchical in structure
- interpretation of patterns depends on context
- essential for complex recognition tasks



Assignment Project Exam Help

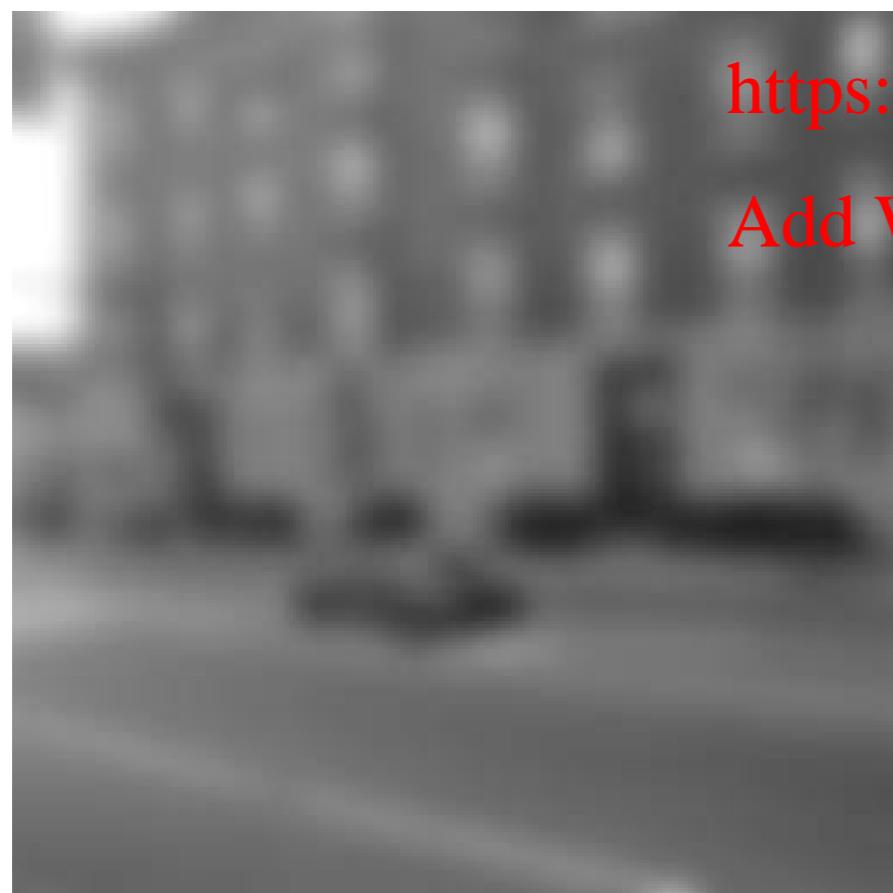
<https://powcoder.com>

Add WeChat powcoder



Motivation: hierarchical, context dependent representations

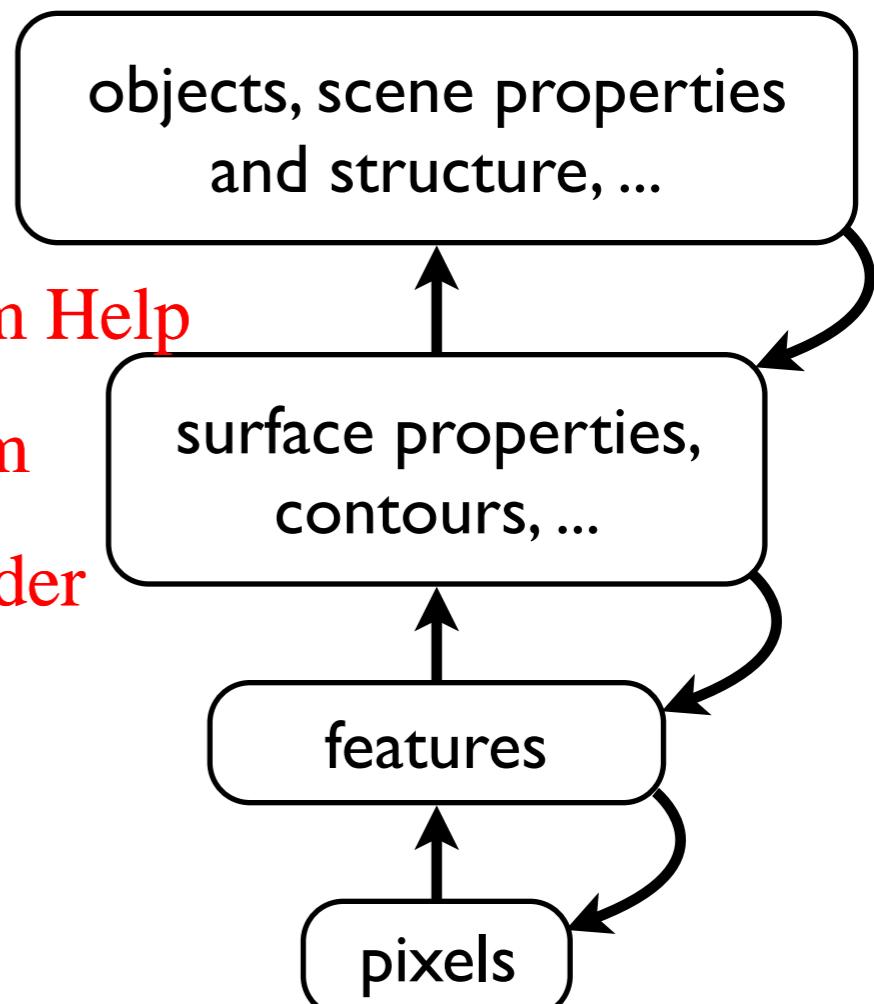
- many real-world patterns are hierarchical in structure
- interpretation of patterns depends on context
- essential for complex recognition tasks



Assignment Project Exam Help

<https://powcoder.com>

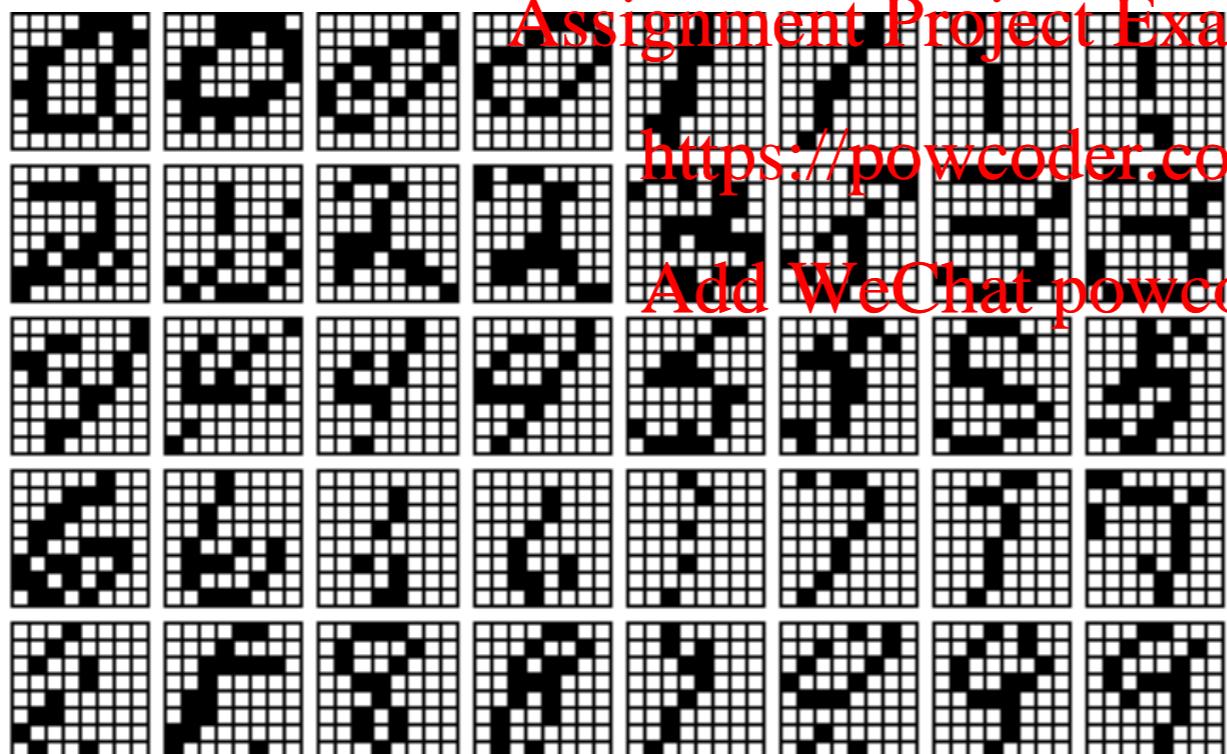
Add WeChat powcoder



Toy problems - Handwritten digits

I - Data: pixels (black or white)

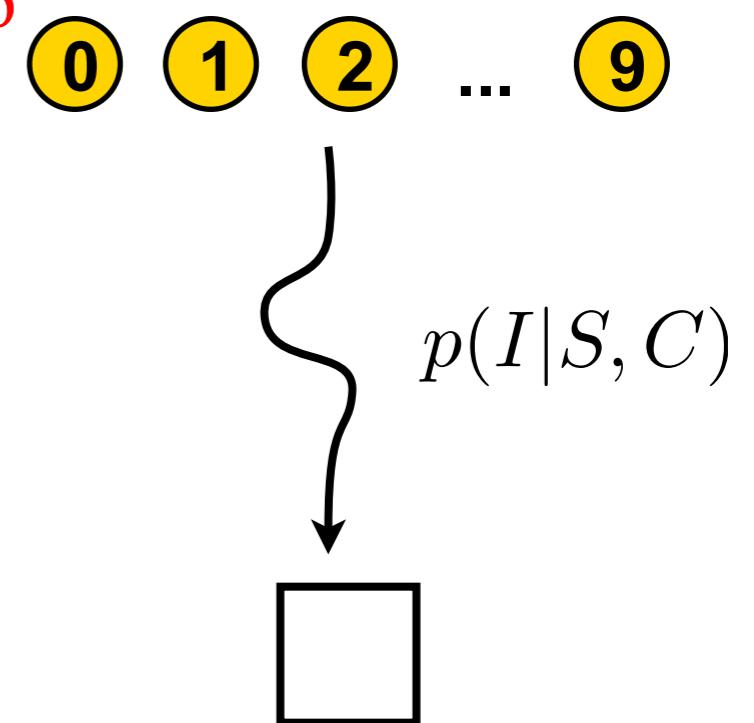
S - Complex “random” patterns with simple underlying structure



Assignment Project Exam Help

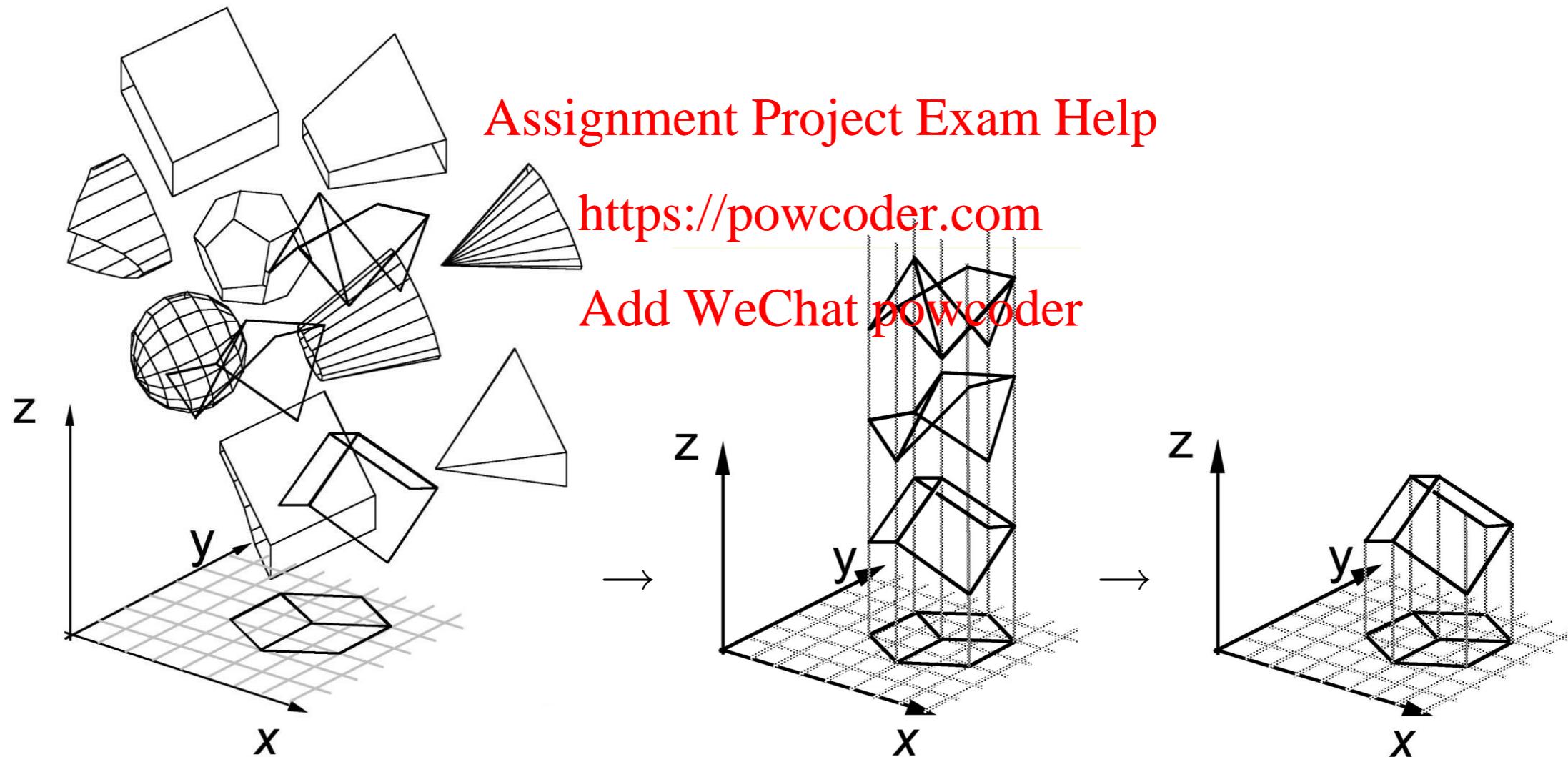
<https://powcoder.com>

Add WeChat powcoder



The inferred explanation is the most probable scene

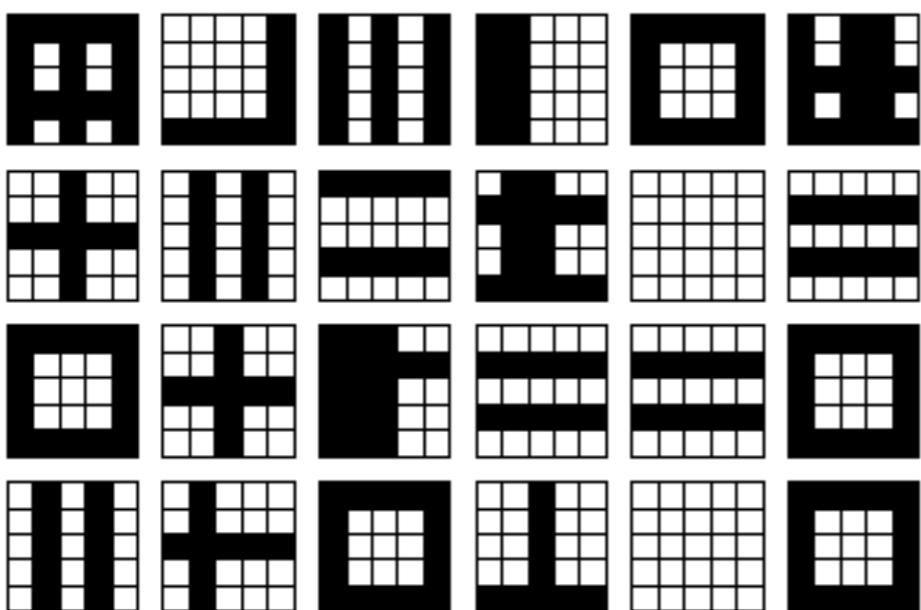
$$p(\hat{S}|I, C) = \arg \max_S \frac{p(I|S, C)p(S|C)}{p(I|C)}$$



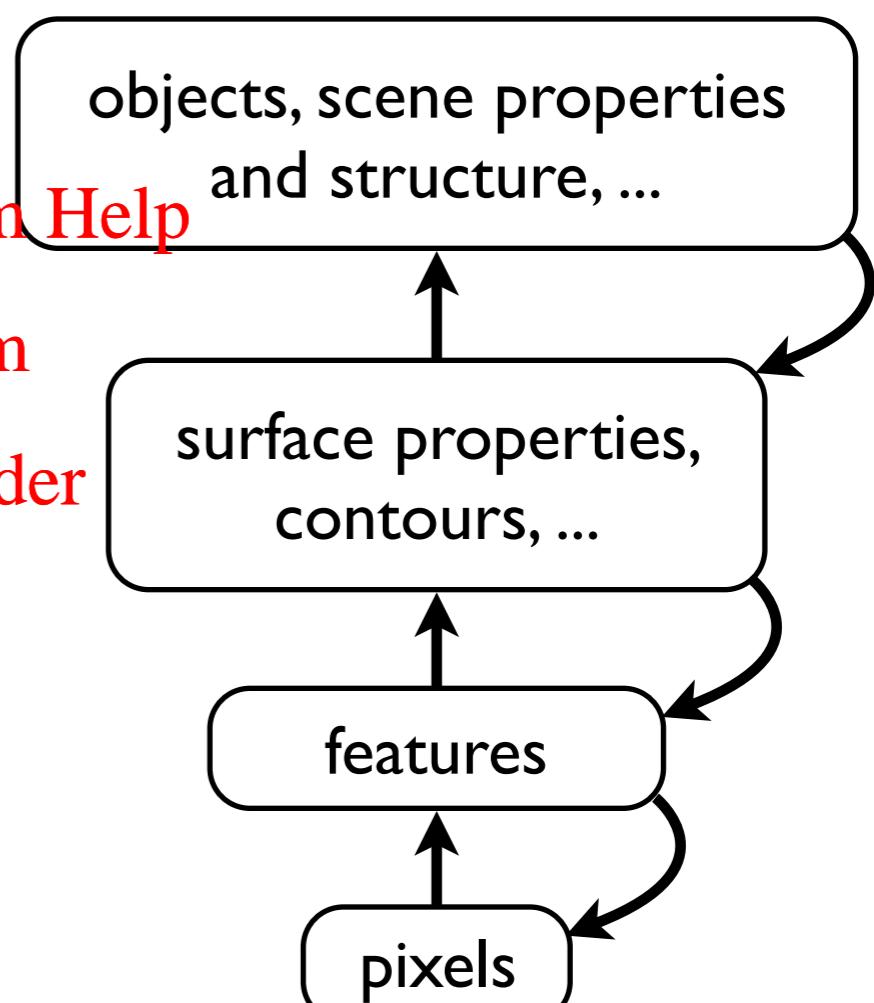
Motivation: learn hierarchical, context dependent representations

- many real-world patterns are hierarchical in structure
- interpretation of patterns depends on context
- essential for complex recognition tasks

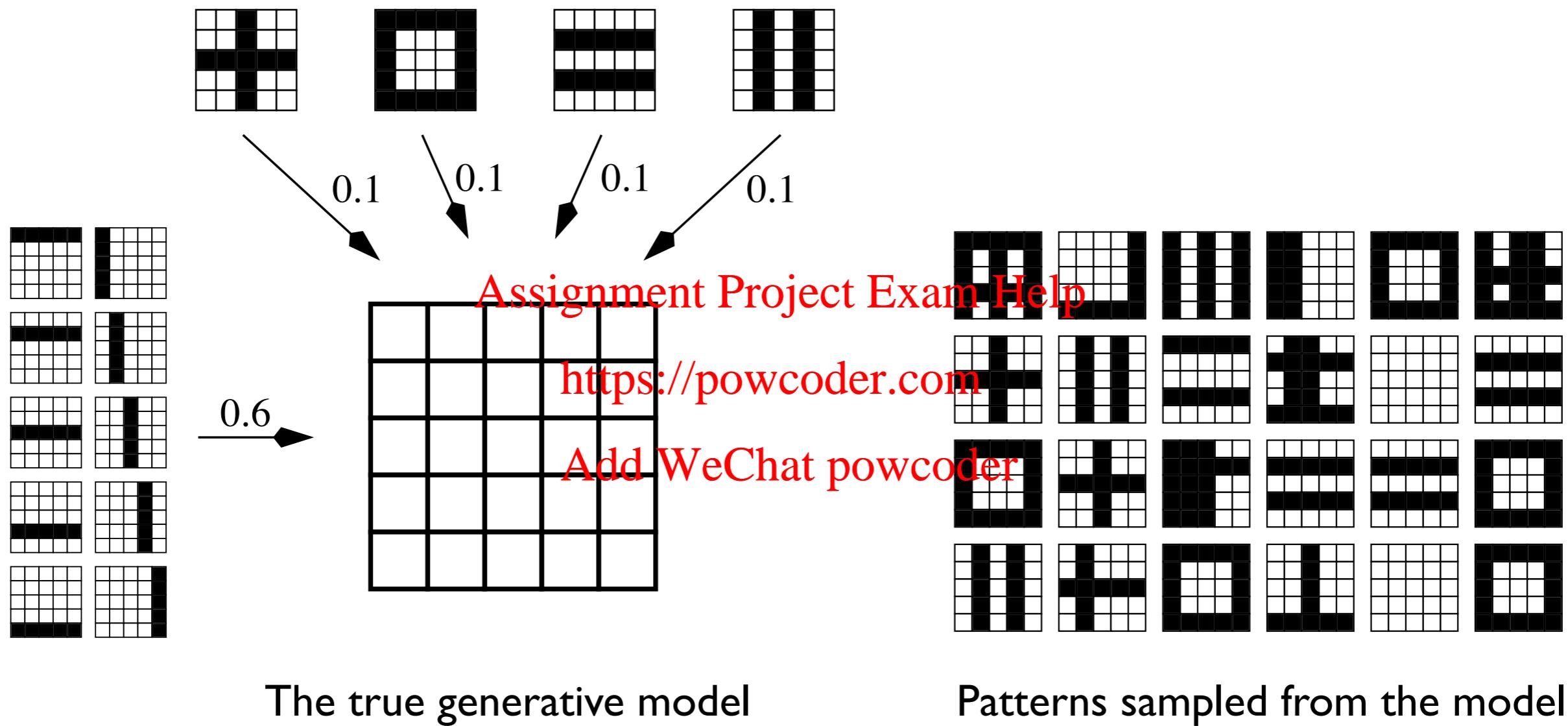
The toy problem:
is it a pattern or a collection of features?



Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

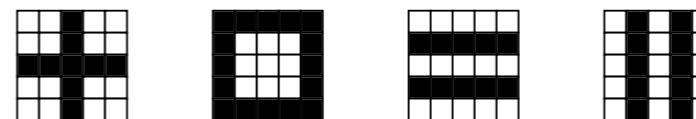


The higher-order lines problem

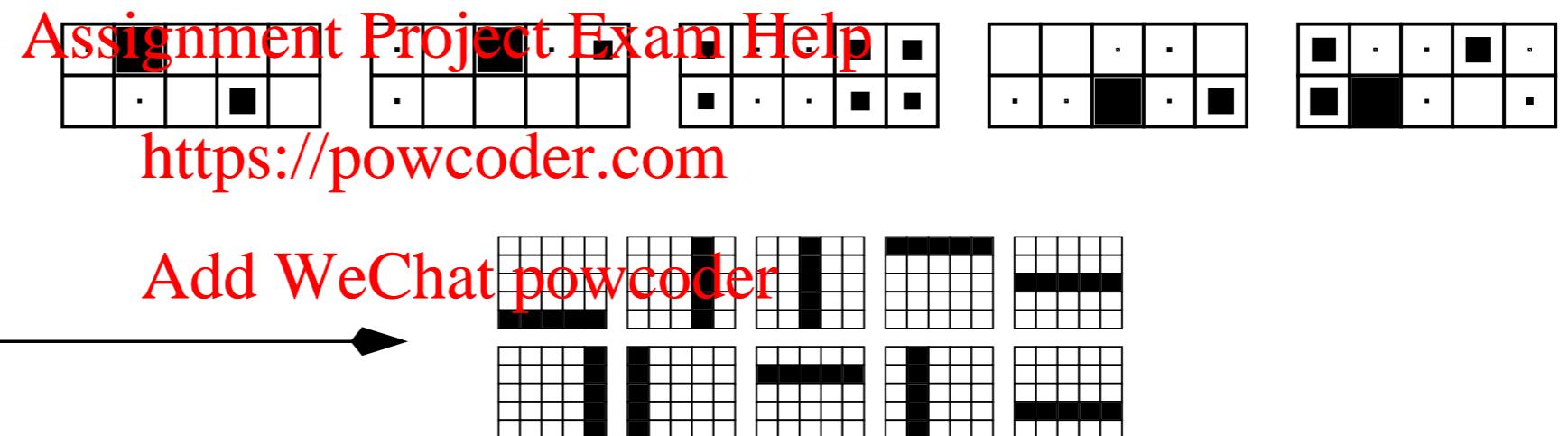
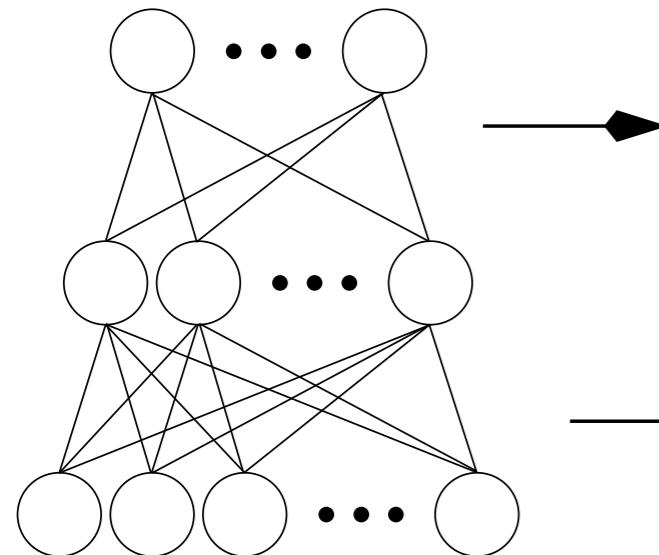


Can we infer the structure of the network given only the patterns?

Weights in a 25-10-5 belief network after learning

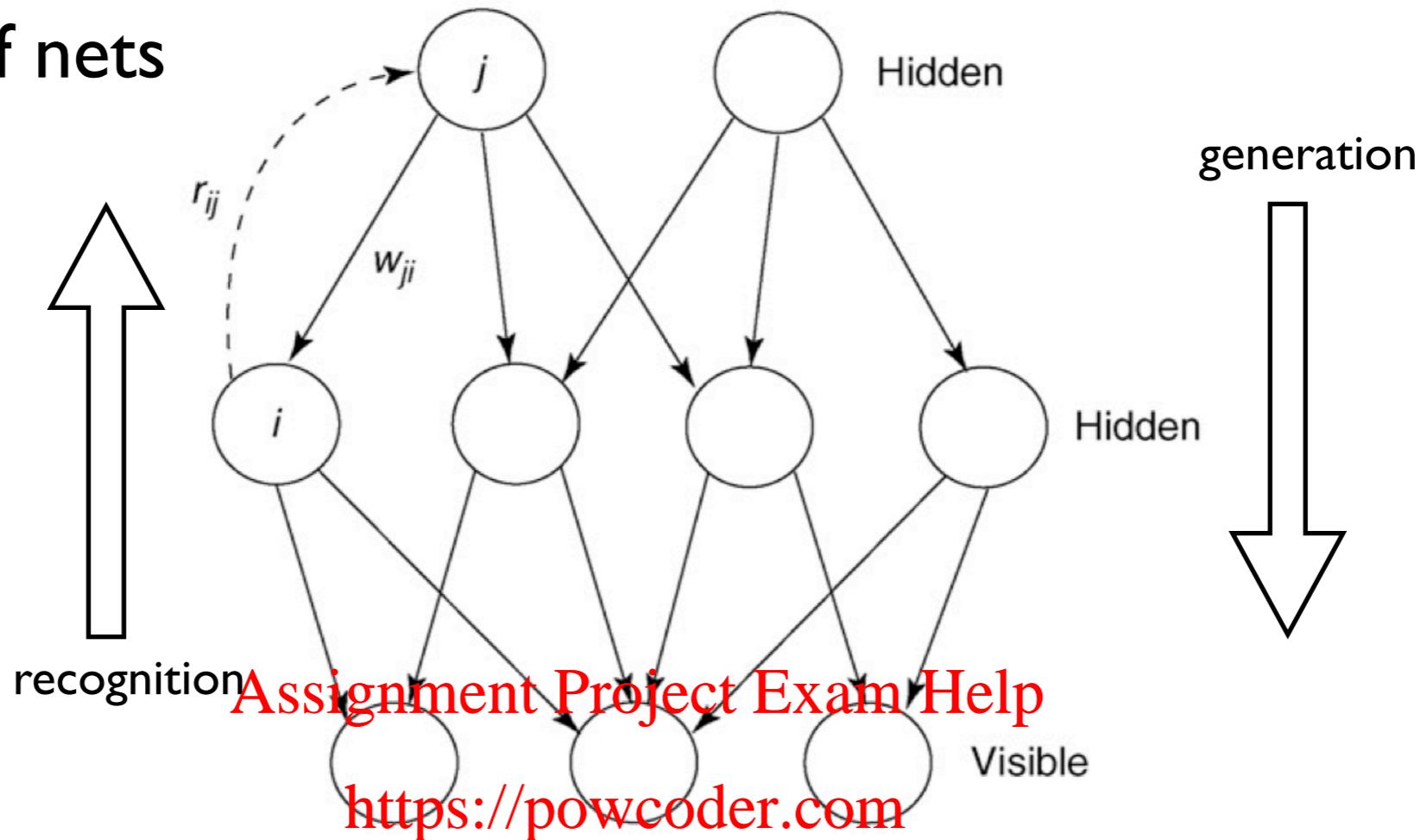


The second layer learns combinations of the first layer features



The first layer of weights learn that patterns are combinations of lines.

Logistic belief nets



- generative model:

$$p(v_i = 1) = \sigma(b_i + \sum_j h_j w_{ij})$$

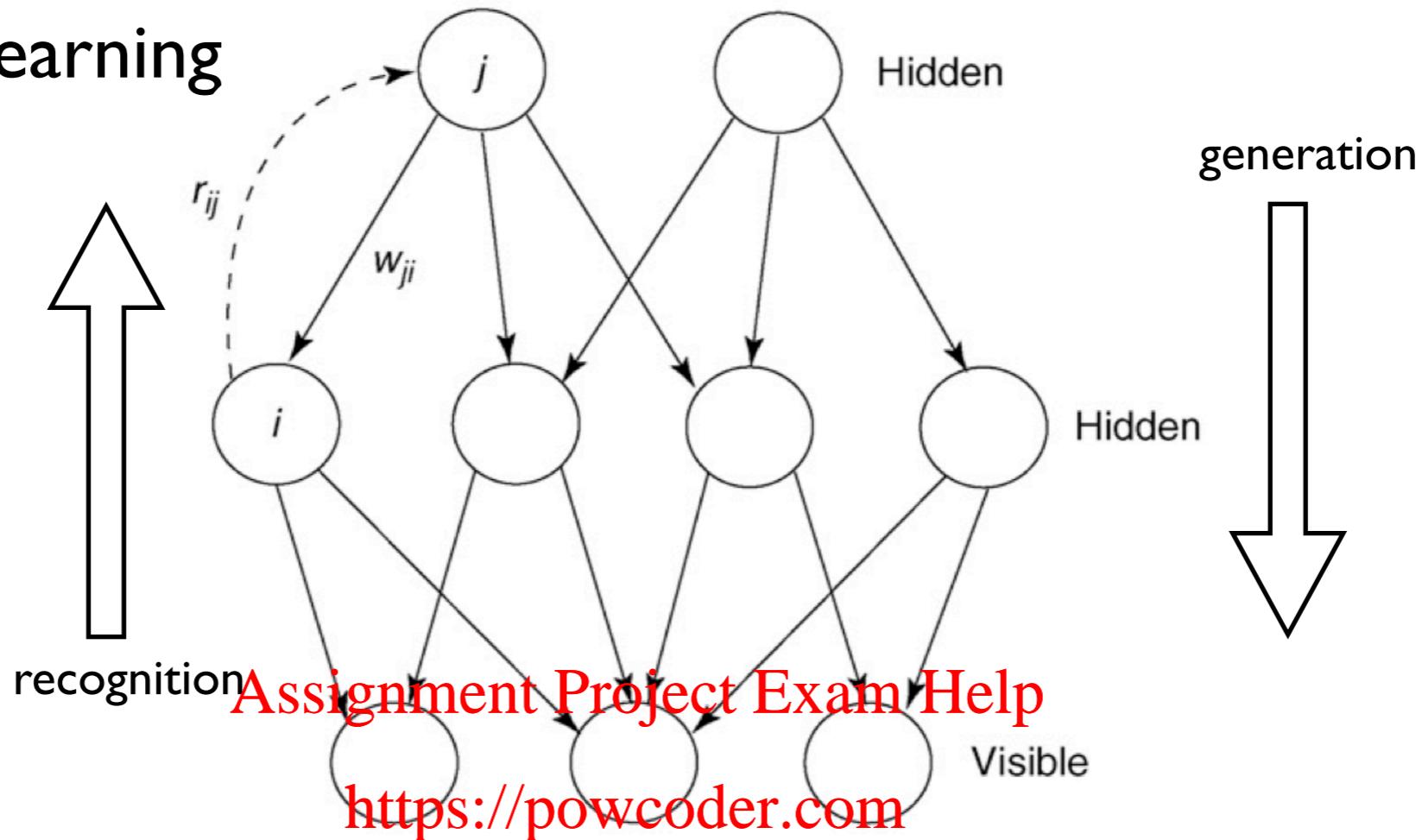
- $\sigma(x)$ is the logistic function:

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

Add WeChat powcoder

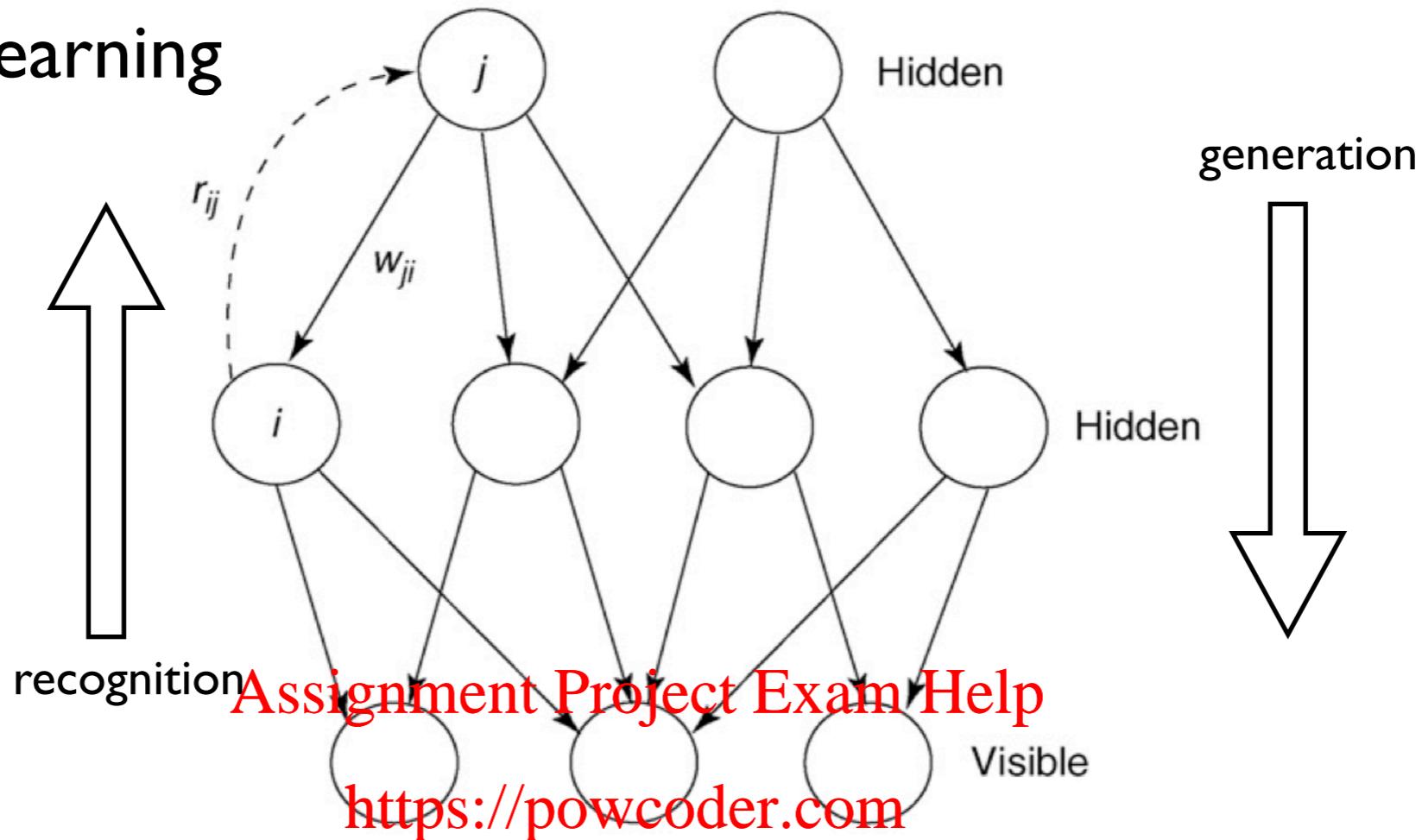
- with deep networks, it is possible to learn complex joint probability distributions

Wake-sleep learning



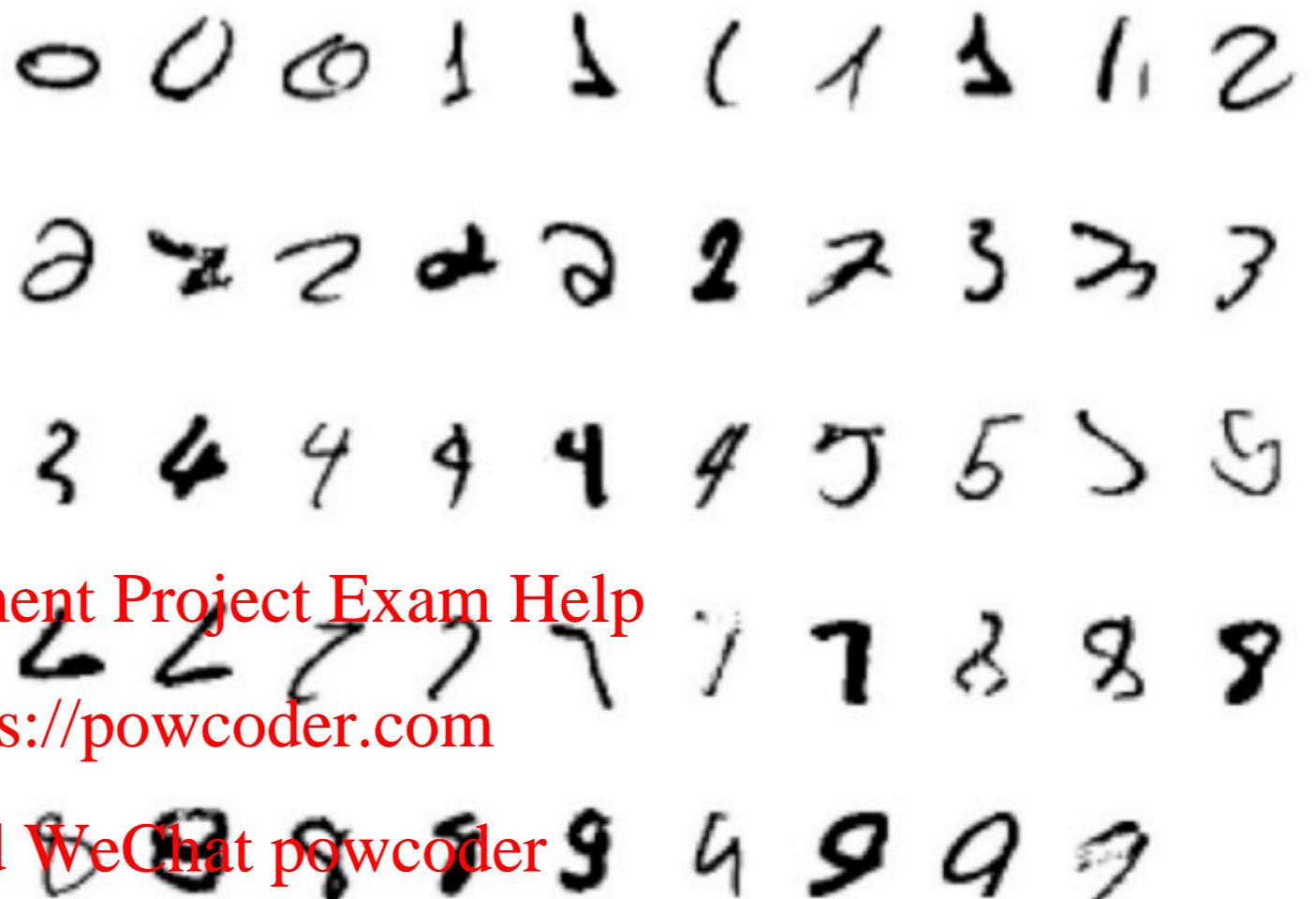
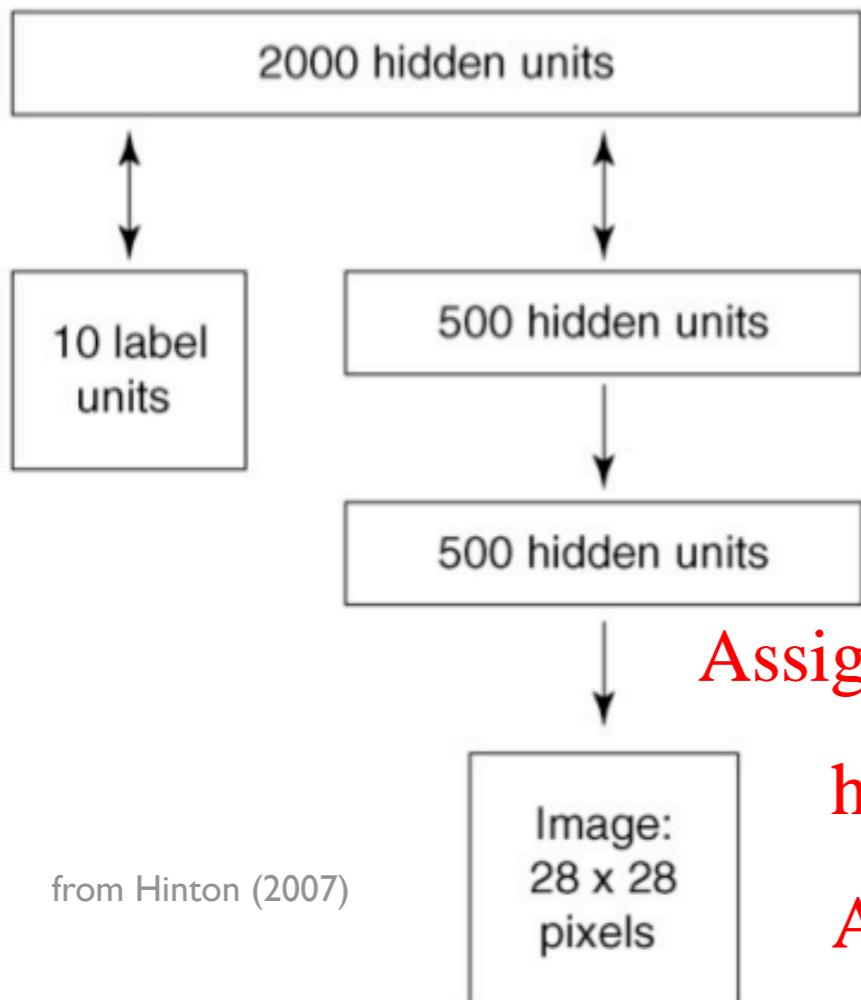
- For probabilistic models [Add WeChat powcoder](#)
 - top-down weights generate patterns from model distribution
 - bottom-up weights convey distribution of data-vectors
 - ideally the two distributions should match
- The “wake-sleep” algorithm adjusts the weights so that the distribution from recognition (wake) matches the distribution from generation (sleep)

Wake-sleep learning



- For each digit in training set [Add WeChat powcoder](#)
 - bottom-up pass: use recognition weights to stochastically set hidden states
$$p(h_j = 1) = \sigma(b_j + \sum_i v_i w_{ij})$$
 - adjust generative weights to improve how model generates training data:
$$\Delta w_{ji} \propto h_j(h_i - \hat{h}_i)$$
 - \hat{h}_i is the probability of activating state i given inferred states h_j

Generative model for hand-written digits



- Generation:

- use alternating Gibbs sampling from top-level assoc. memory
- use directed weights to stochastically generate pixel probs. from sampled binary of 500 hidden units

- Recognition:

- Use bottom-up weights to produce binary activities in two lower layers
- use alternating Gibbs sampling in the top two layers

Demo: deep-belief network model (Hinton)



► ▶ INCREASE SPEED DETAILED VIEW

