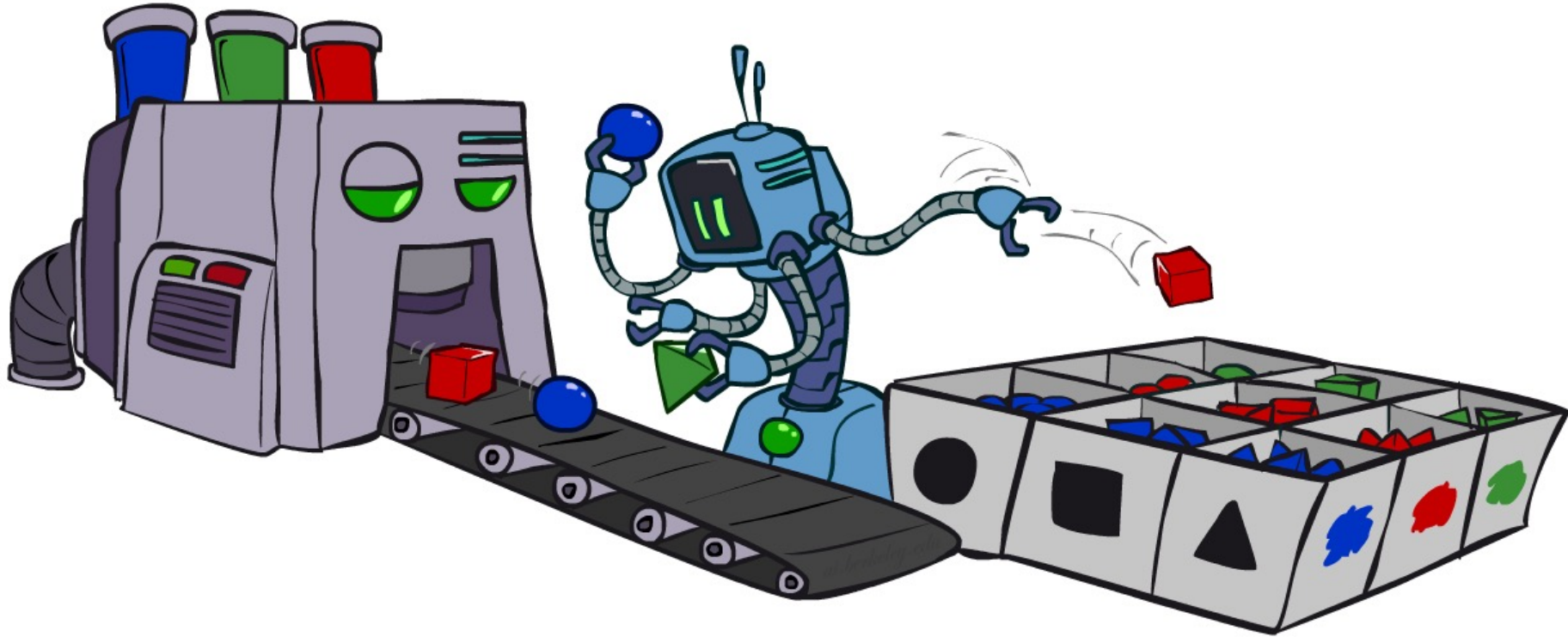# CS 188: Artificial Intelligence

# Bayes' Nets: Sampling

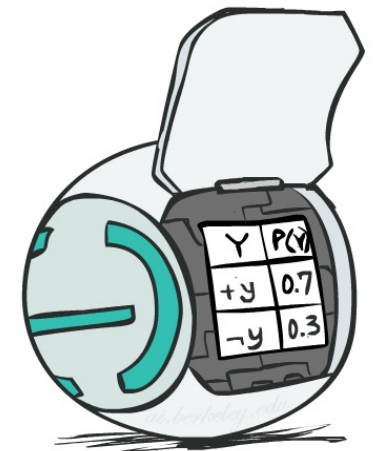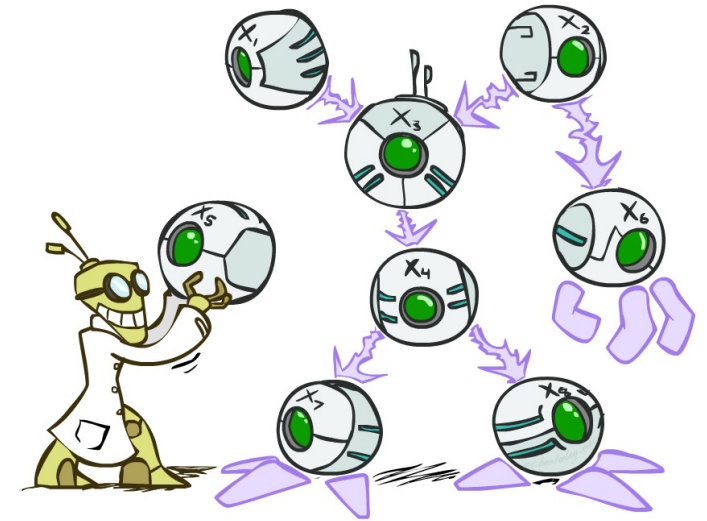Instructor: Professor Dragan --- University of California, Berkeley

# Bayes' Net Representation

- A directed, acyclic graph, one node per random variable

- A conditional probability table (CPT) for each node

  - A collection of distributions over X, one for each combination of parents' values
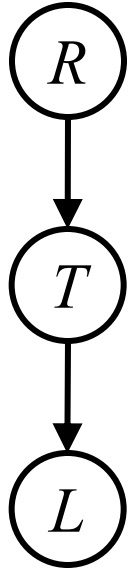  $$P(X|a_1 \ldots a_n)$$

- Bayes' nets implicitly encode joint distributions

  - As a product of local conditional distributions

  - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

  $$P(x_1, x_2, \ldots x_n) = \prod_{i=1}^{n} P(x_i | parents(X_i))$$

# Recap: Bayesian Inference (Exact)

$R$

$T$

$L$

$$P(L) = \ ?$$

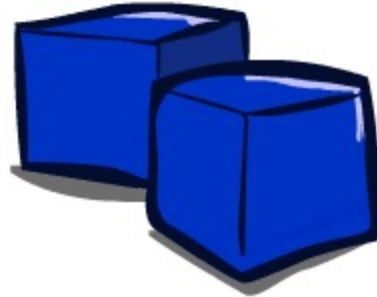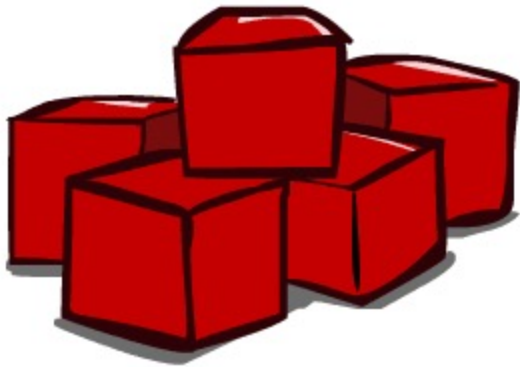- Inference by Enumeration

$$= \sum_t \sum_r P(L|t)P(r)P(t|r)$$

Join on r

Join on t

Eliminate r

Eliminate t

- Variable Elimination

$$= \sum_t P(L|t) \sum_r P(r)P(t|r)$$

Join on r

Eliminate r

Join on t

Eliminate t

# Approximate Inference: Sampling

# Sampling

- Sampling is a lot like repeated simulation

  - Predicting the weather, basketball games, …

- Basic idea

  - Draw N samples from a sampling distribution S

  - Compute an approximate posterior probability

  - Show this converges to the true probability P

- Why sample?

  - Learning: get samples from a distribution you don't know

  - Inference: getting a sample is faster than computing the right answer (e.g. with variable elimination)

# Sampling

- Sampling from given distribution

  - Step 1: Get sample $u$ from uniform distribution over $[0, 1)$
    - E.g. random() in python

  - Step 2: Convert this sample $u$ into an outcome for the given distribution by having each target outcome associated with a sub-interval of $[0,1)$ with sub-interval size equal to probability of the outcome

- Example

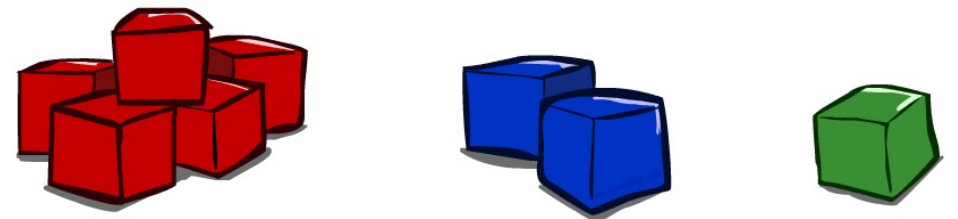| C | P(C) |
|-------|------|
| red | 0.6 |
| green | 0.1 |
| blue | 0.3 |

$$0 \leq u < 0.6, \rightarrow C = red$$
$$0.6 \leq u < 0.7, \rightarrow C = green$$
$$0.7 \leq u < 1, \rightarrow C = blue$$

- If random() returns $u = 0.83$, then our sample is $C = blue$
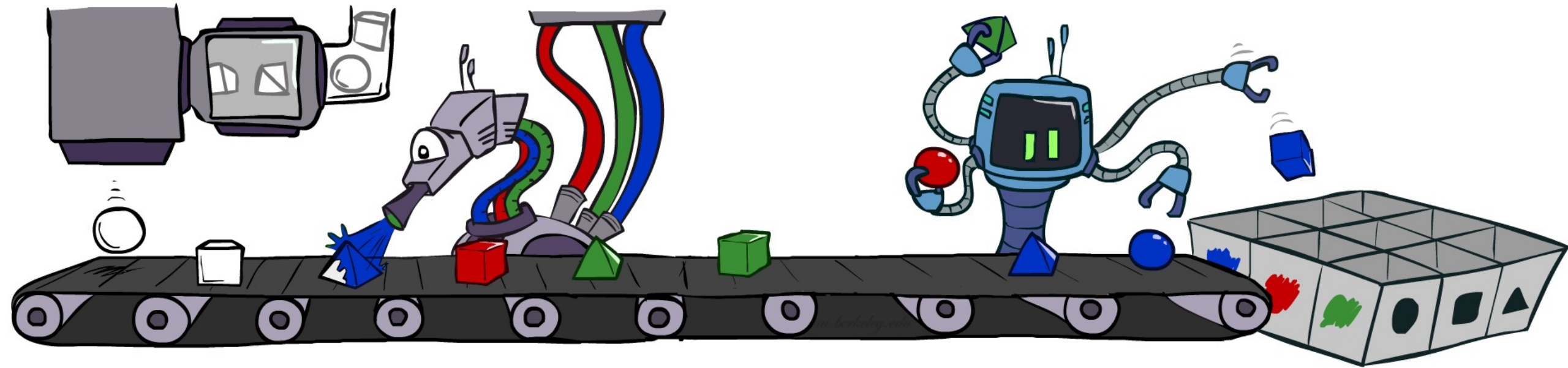- E.g, after sampling 8 times:

# Sampling in Bayes' Nets

- Prior Sampling

- Rejection Sampling

- Likelihood Weighting

- Gibbs Sampling

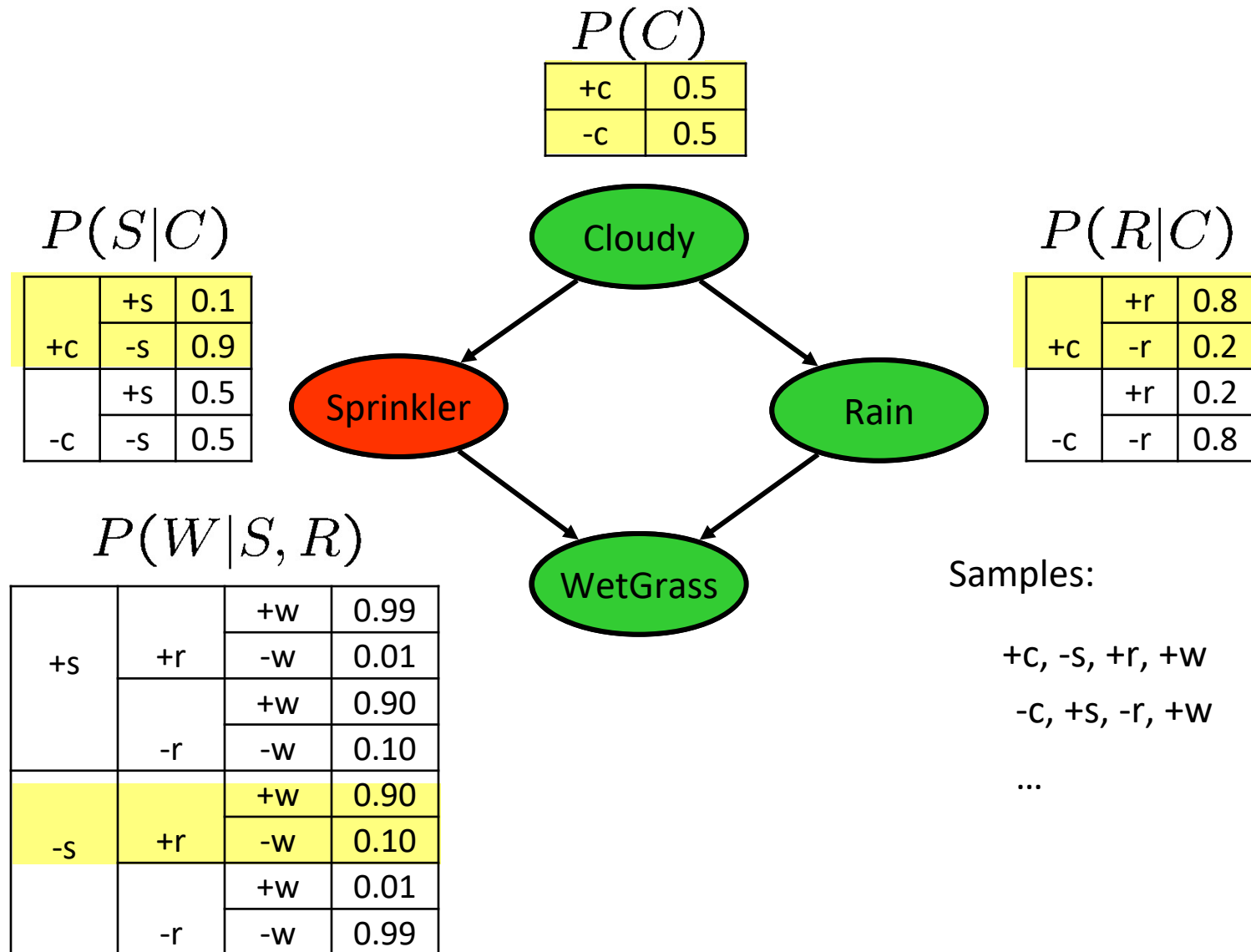# Prior Sampling

# Prior Sampling

$P(C)$

| +c | 0.5 |
|----|-----|
| -c | 0.5 |

$P(S|C)$

| | +s | 0.1 |
|----|----|-----|
| +c | -s | 0.9 |
| | +s | 0.5 |
| -c | -s | 0.5 |

Cloudy

Sprinkler

Rain

WetGrass

$P(R|C)$

| | +r | 0.8 |
|----|----|-----|
| +c | -r | 0.2 |
| | +r | 0.2 |
| -c | -r | 0.8 |

$P(W|S, R)$

| | | +w | 0.99 |
|----|----|----|------|
| +s | +r | -w | 0.01 |
| | | +w | 0.90 |
| | -r | -w | 0.10 |
| | | +w | 0.90 |
| -s | +r | -w | 0.10 |
| | | +w | 0.01 |
| | -r | -w | 0.99 |

Samples:

+c, -s, +r, +w

-c, +s, -r, +w
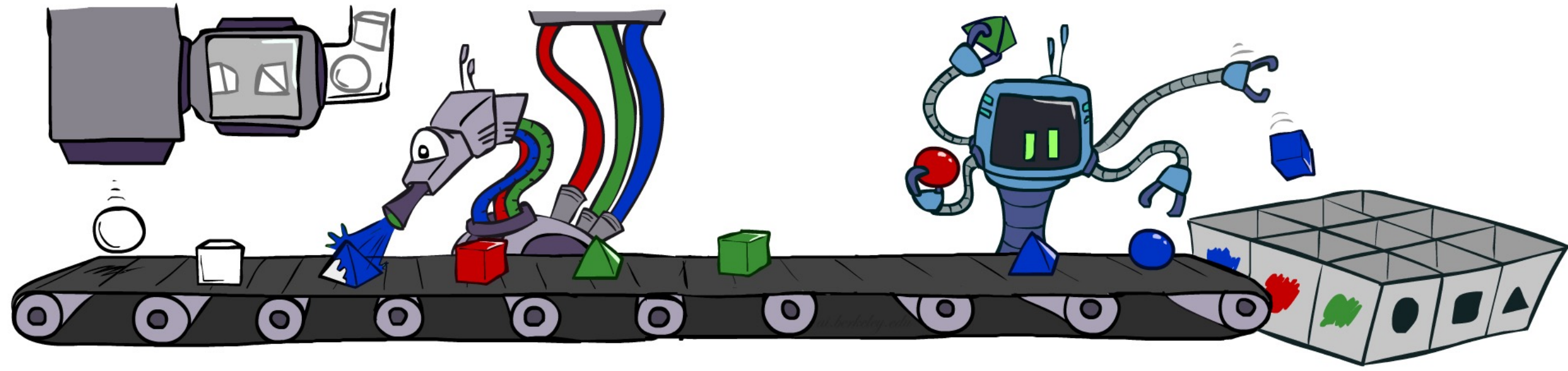
…

# Prior Sampling

- For i = 1, 2, …, n in topological order

  - Sample $x_i$ from $P(X_i \mid \text{Parents}(X_i))$

- Return $(x_1, x_2, …, x_n)$

# Prior Sampling

- This process generates samples with probability:

$$S_{PS}(x_1 \ldots x_n) = \prod_{i=1}^{n} P(x_i | \mathsf{Parents}(X_i)) = P(x_1 \ldots x_n)$$

  …i.e. the BN's joint probability

- Let the number of samples of an event be $N_{PS}(x_1 \ldots x_n)$

- Then
$$\lim_{N \to \infty} \widehat{P}(x_1, \ldots, x_n) = \lim_{N \to \infty} N_{PS}(x_1, \ldots, x_n)/N$$
$$= S_{PS}(x_1, \ldots, x_n)$$
$$= P(x_1 \ldots x_n)$$

- I.e., the sampling procedure is consistent

# Example

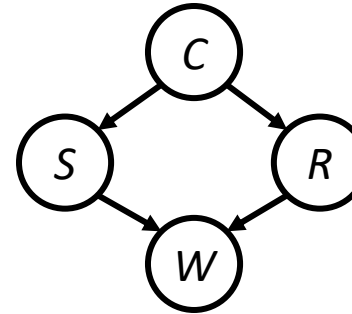- We'll get a bunch of samples from the BN:

  +c, -s, +r, +w
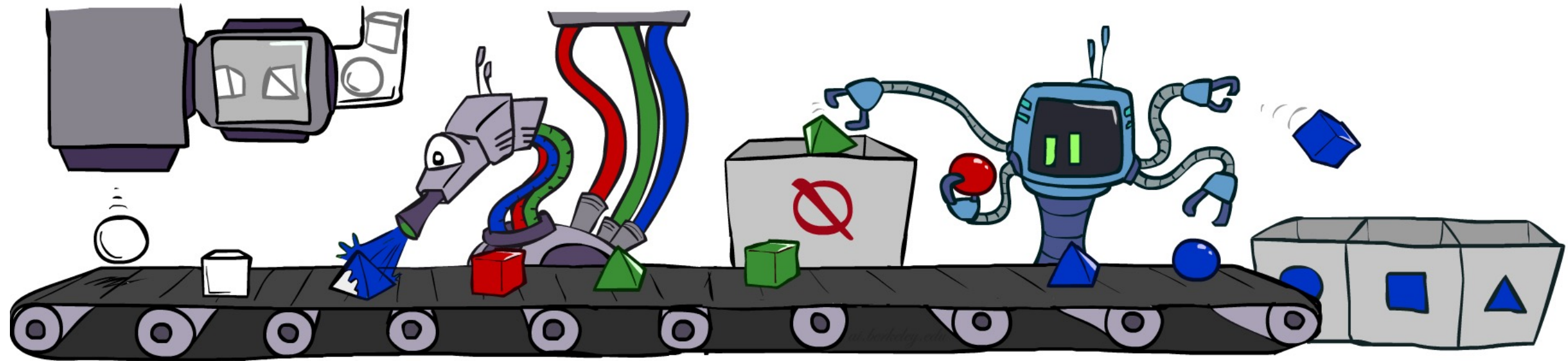
  +c, +s, +r, +w

  -c, +s, +r, -w

  +c, -s, +r, +w

  -c, -s, -r, +w

- If we want to know P(W)
  - We have counts <+w:4, -w:1>
  - Normalize to get P(W) = <+w:0.8, -w:0.2>
  - This will get closer to the true distribution with more samples
  - Can estimate anything else, too
    - P(C | +w)?   P(C | +r, +w)?
    - Can also use this to estimate expected value of f(X) - Monte Carlo Estimation
  - What about P(C | -r, -w)?

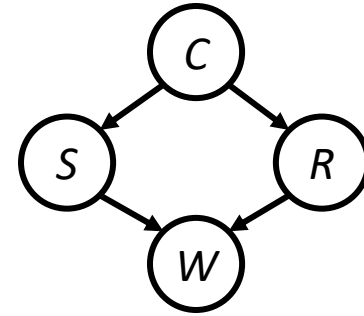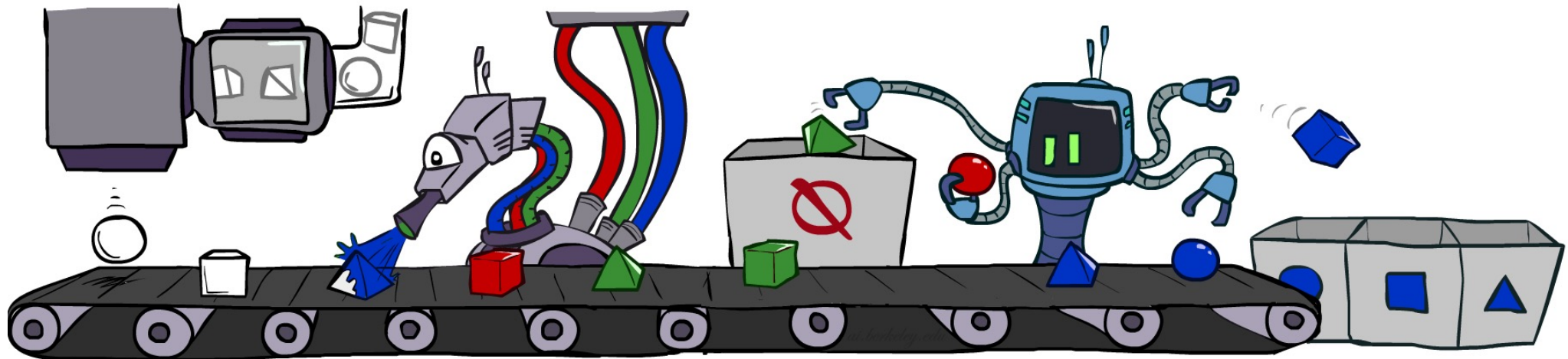# Rejection Sampling

- Let's say we want P(C)
  - Just tally counts of C as we go

- Let's say we want P(C | +s)
  - Same thing: tally C outcomes, but ignore (reject) samples which don't have S=+s
  - This is called rejection sampling
  - We can toss out samples early!
  - It is also consistent for conditional probabilities (i.e., correct in the limit)
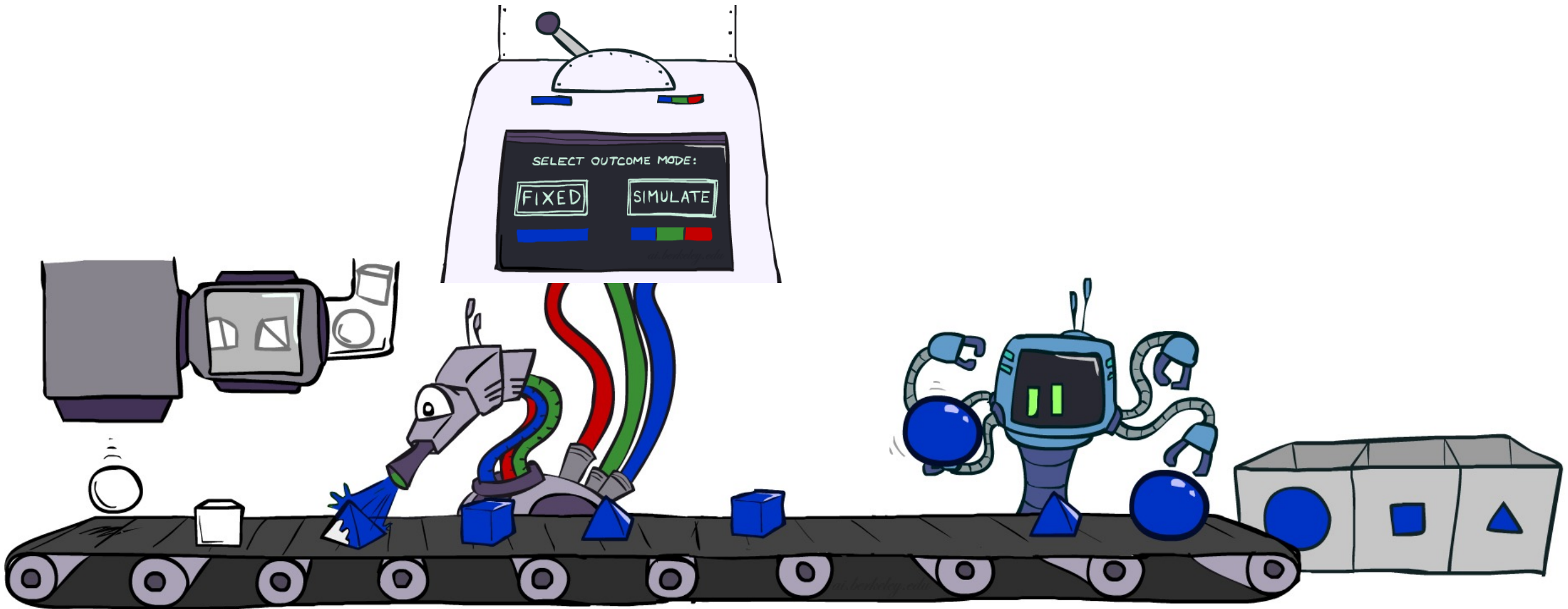
# Rejection Sampling

- Input: evidence instantiation
- For i = 1, 2, …, n in topological order

  - Sample $x_i$ from $P(X_i \mid Parents(X_i))$
  - If $x_i$ not consistent with evidence
    - Reject: return – no sample is generated in this cycle
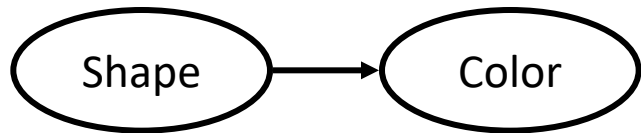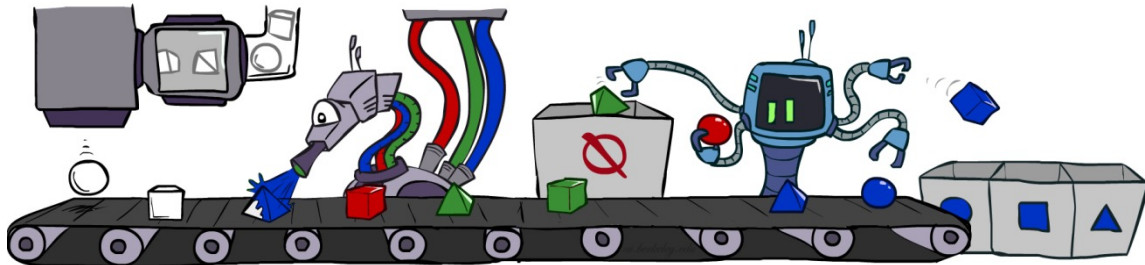
- Return $(x_1, x_2, …, x_n)$

# Likelihood Weighting



SELECT OUTCOME MODE:

FIXED    SIMULATE

# Likelihood Weighting

- ## Problem with rejection sampling:
  - ### If evidence is unlikely, rejects lots of samples
  - ### Consider P( Shape | blue )

- ## Idea: fix evidence variables and sample the rest
  - ### Problem: sample distribution not consistent!
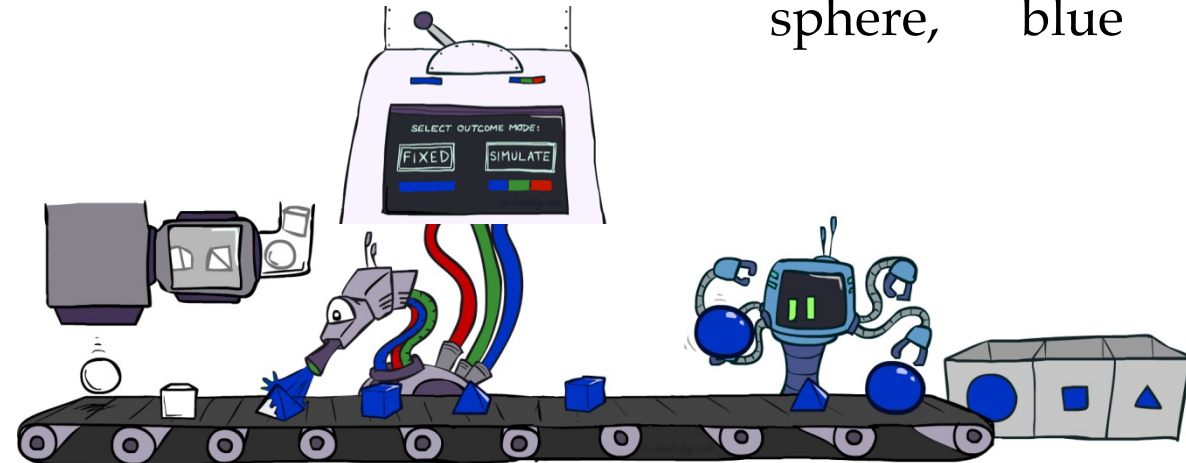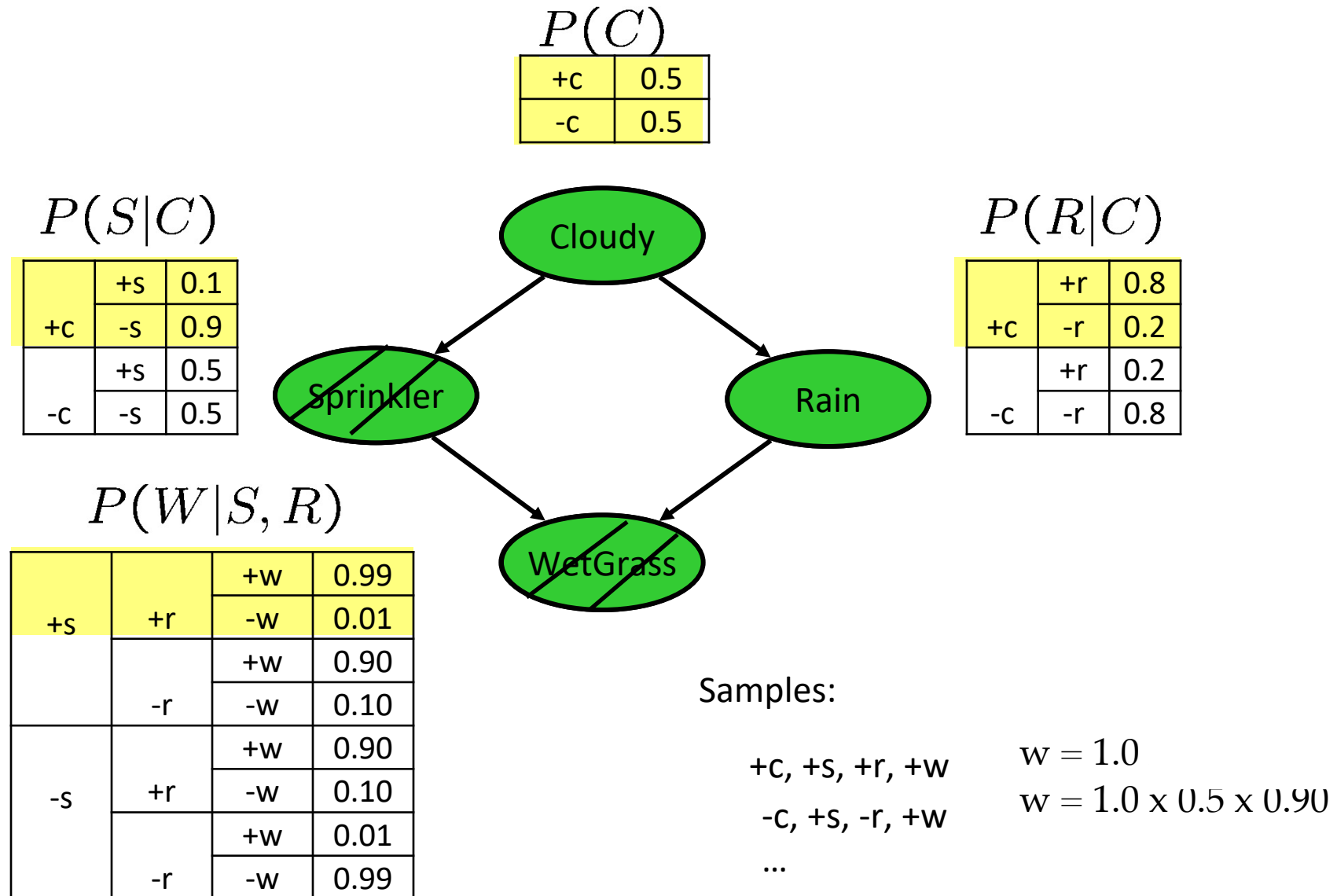  - ### Solution: weight by probability of evidence given parents



~~pyramid, green~~
~~pyramid, red~~
sphere, blue
~~cube, red~~
~~sphere, green~~

pyramid, blue
pyramid, blue
sphere, blue
cube, blue
sphere, blue

# Likelihood Weighting

$P(C)$

| +c | 0.5 |
|----|-----|
| -c | 0.5 |

$P(S|C)$

|    | +s | 0.1 |
|----|----|-----|
| +c | -s | 0.9 |
|    | +s | 0.5 |
| -c | -s | 0.5 |

$P(R|C)$

|    | +r | 0.8 |
|----|----|-----|
| +c | -r | 0.2 |
|    | +r | 0.2 |
| -c | -r | 0.8 |

Cloudy

Sprinkler

Rain

WetGrass

$P(W|S,R)$

|    |    | +w | 0.99 |
|----|----|----|------|
| +s | +r | -w | 0.01 |
|    |    | +w | 0.90 |
|    | -r | -w | 0.10 |
|    |    | +w | 0.90 |
| -s | +r | -w | 0.10 |
|    |    | +w | 0.01 |
|    | -r | -w | 0.99 |

Samples:

+c, +s, +r, +w

-c, +s, -r, +w

...

$w = 1.0$

$w = 1.0 \times 0.5 \times 0.90$
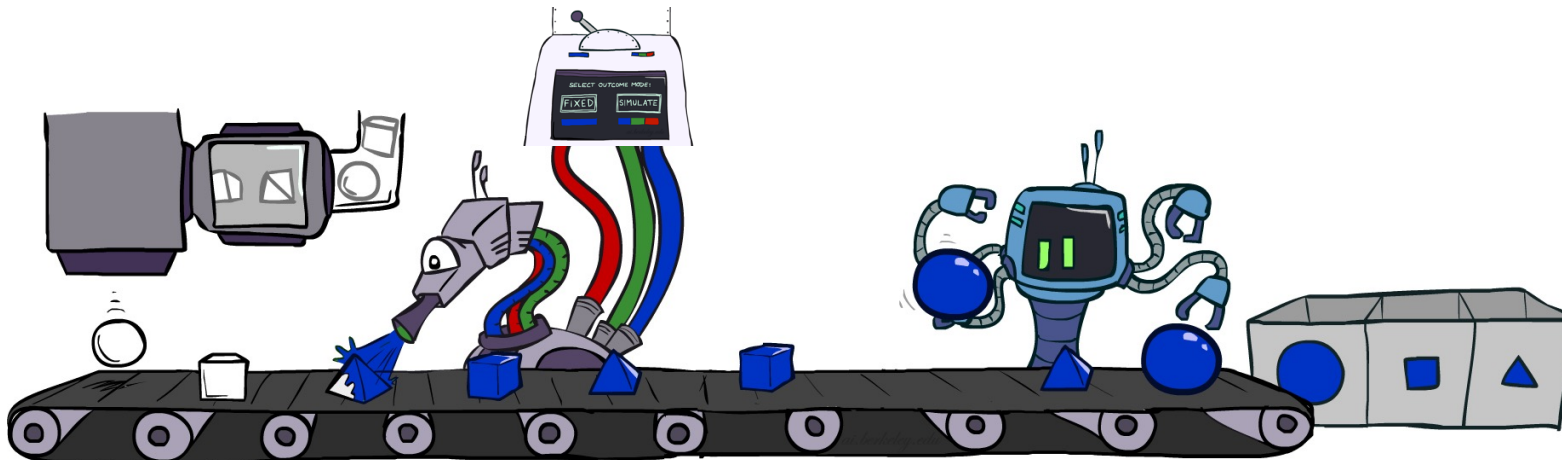
# Likelihood Weighting

- Input: evidence instantiation
- $w = 1.0$
- for $i = 1, 2, \ldots, n$ in topological order
  - if $X_i$ is an evidence variable
    - $X_i$ = observation $x_i$ for $X_i$
    - Set $w = w * P(x_i \mid \text{Parents}(X_i))$
  - else
    - Sample $x_i$ from $P(X_i \mid \text{Parents}(X_i))$
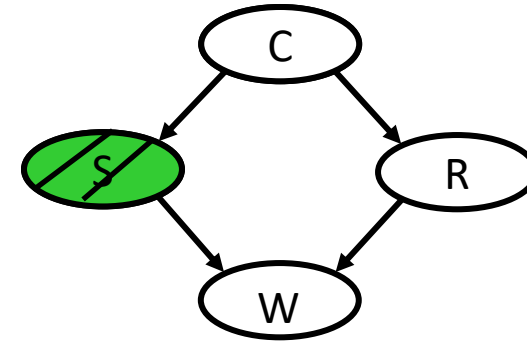- return $(x_1, x_2, \ldots, x_n)$, $w$

# Likelihood Weighting

- Sampling distribution if z sampled and e fixed evidence

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{l} P(z_i | \text{Parents}(Z_i))$$

- Now, samples have weights

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{m} P(e_i | \text{Parents}(E_i))$$

- Together, weighted sampling distribution is consistent

$$S_{\text{WS}}(z, e) \cdot w(z, e) = \prod_{i=1}^{l} P(z_i | \text{Parents}(z_i)) \prod_{i=1}^{m} P(e_i | \text{Parents}(e_i))$$
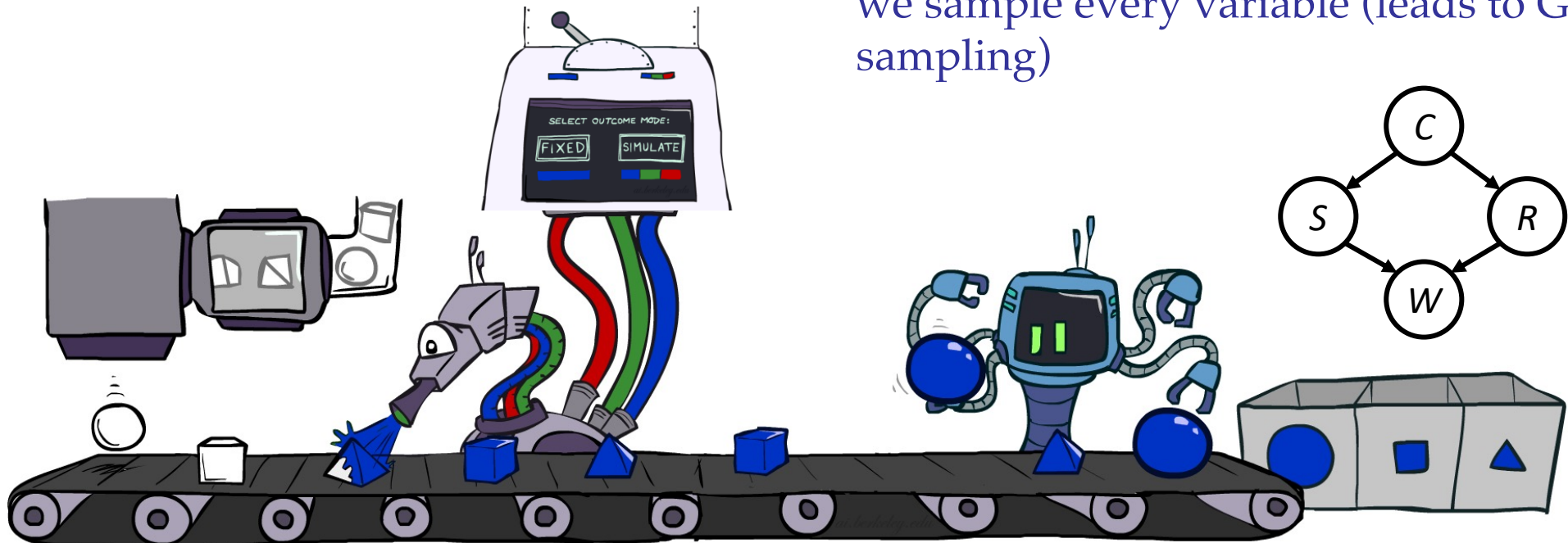
$$= P(\mathbf{z}, \mathbf{e})$$

# Likelihood Weighting

- Likelihood weighting is helpful
  - We have taken evidence into account as we generate the sample
  - E.g. here, W's value will get picked based on the evidence values of S, R
  - More of our samples will reflect the state of the world suggested by the evidence
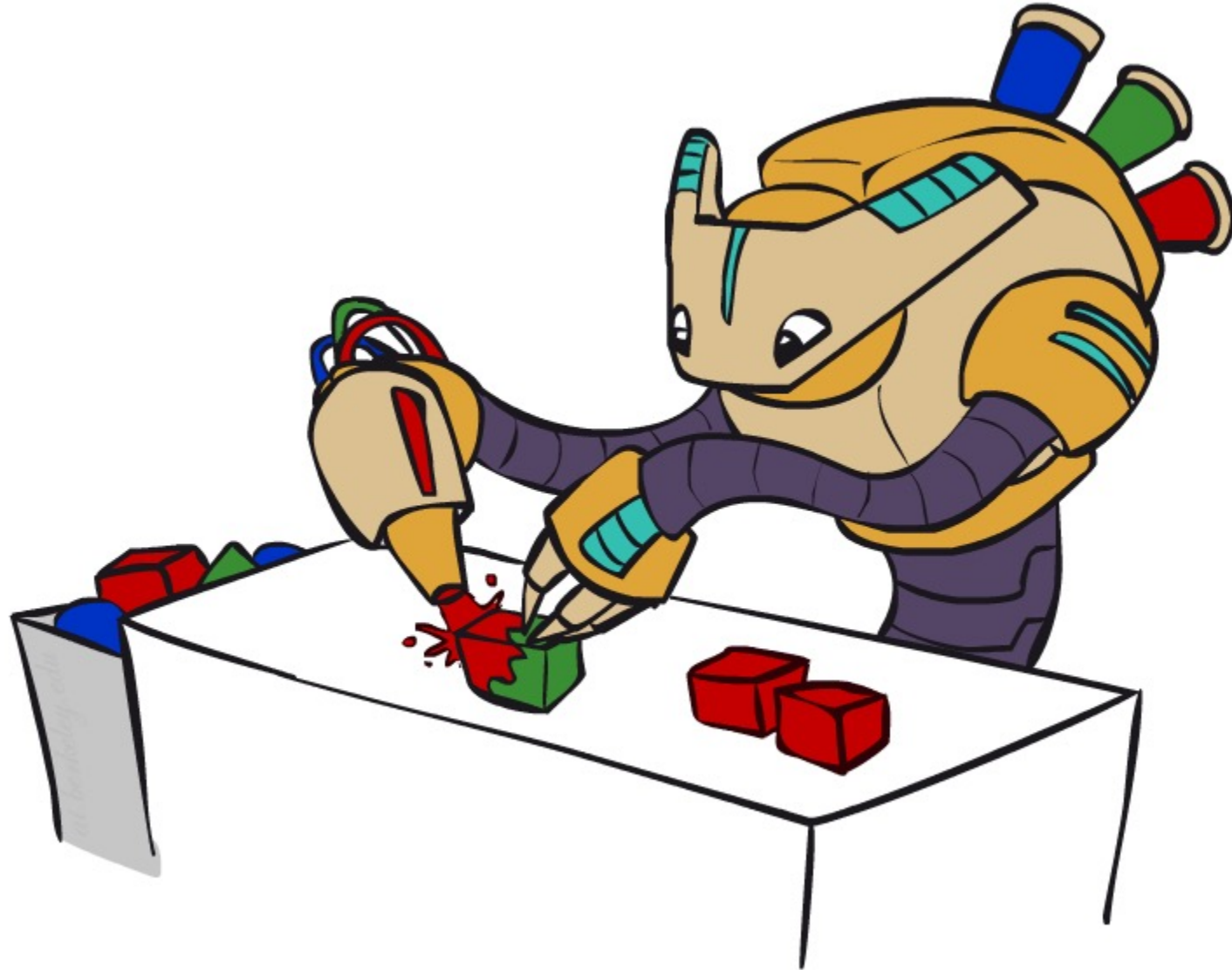
- Likelihood weighting doesn't solve all our problems
  - Evidence influences the choice of downstream variables, but not upstream ones (C isn't more likely to get a value matching the evidence)
- We would like to consider evidence when we sample every variable (leads to Gibbs sampling)
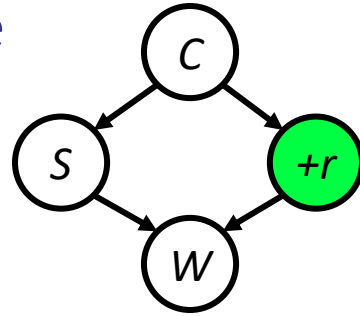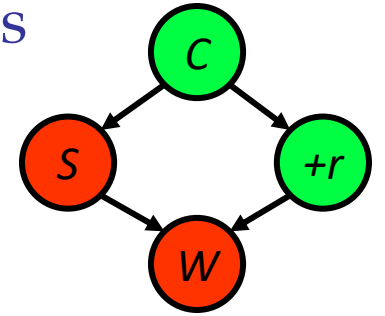
# Gibbs Sampling

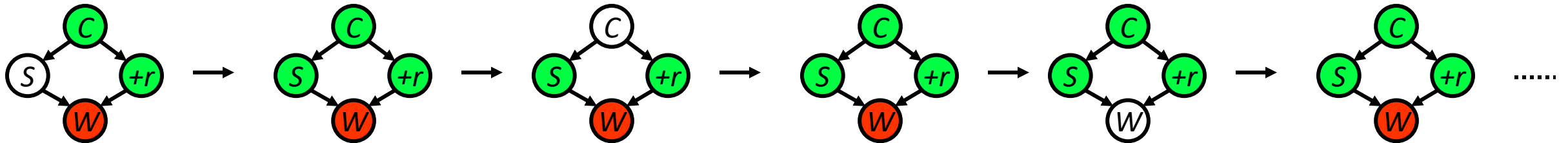# Gibbs Sampling Example: P( S | +r)

- **Step 1: Fix evidence**
  - R = +r



- **Step 2: Initialize other variables**
  - Randomly



- **Steps 3: Repeat**
  - Choose a non-evidence variable X
  - Resample X from P( X | all other variables)*

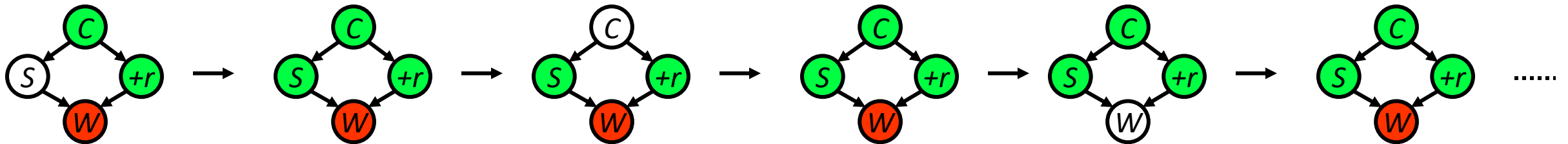

Sample from $P(S|+c,-w,+r)$   Sample from $P(C|+s,-w,+r)$   Sample from $P(W|+s,+c,+r)$
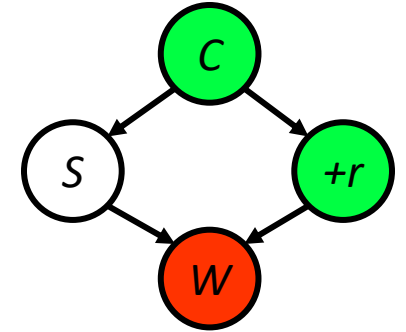
# Gibbs Sampling

- **Procedure**: keep track of a full instantiation $x_1, x_2, \ldots, x_n$. Start with an arbitrary instantiation consistent with the evidence. Sample one variable at a time, conditioned on all the rest, but keep evidence fixed. Keep repeating this for a long time.

- **Property**: in the limit of repeating this infinitely many times the resulting samples come from the correct distribution (i.e. conditioned on evidence).

- **Rationale**: both upstream and downstream variables condition on evidence.

- In contrast: likelihood weighting only conditions on upstream evidence, and hence weights obtained in likelihood weighting can sometimes be very small. Sum of weights over all samples is indicative of how many "effective" samples were obtained, so we want high weight.

# Resampling of One Variable

- Sample from P(S | +c, +r, -w)

$$P(S|+c,+r,-w) = \frac{P(S,+c,+r,-w)}{P(+c,+r,-w)}$$

$$= \frac{P(S,+c,+r,-w)}{\sum_s P(s,+c,+r,-w)}$$

$$= \frac{P(+c)P(S|+c)P(+r|+c)P(-w|S,+r)}{\sum_s P(+c)P(s|+c)P(+r|+c)P(-w|s,+r)}$$

$$= \frac{P(+c)P(S|+c)P(+r|+c)P(-w|S,+r)}{P(+c)P(+r|+c)\sum_s P(s|+c)P(-w|s,+r)}$$

$$= \frac{P(S|+c)P(-w|S,+r)}{\sum_s P(s|+c)P(-w|s,+r)}$$

- Many things cancel out – only CPTs with S remain!
- More generally: only CPTs that have resampled variable need to be considered, and joined together
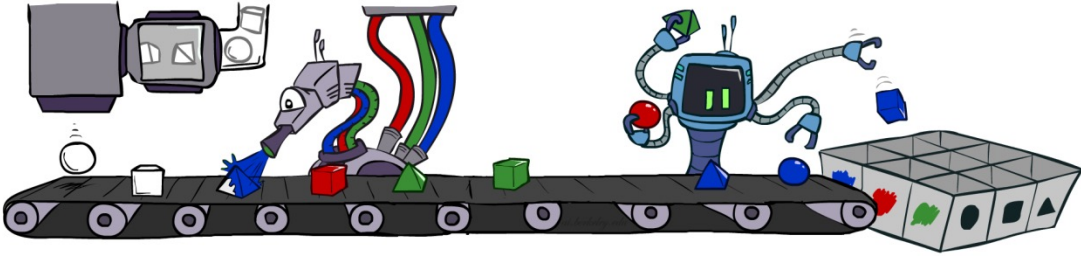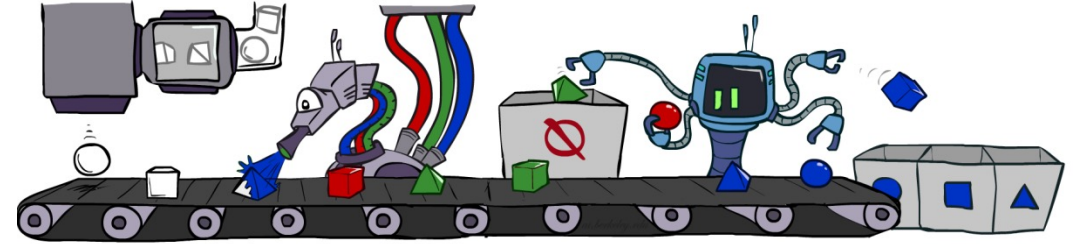
# More Details on Gibbs Sampling*

- Gibbs sampling belongs to a family of sampling methods called Markov chain Monte Carlo (MCMC)

  - Specifically, it is a special case of a subset of MCMC methods called Metropolis-Hastings

- You can read more about this here:

  - https://ermongroup.github.io/cs228-notes/inference/sampling/

# Bayes' Net Sampling Summary

- Prior Sampling  P( Q )



- Likelihood Weighting  P( Q | e)



- Rejection Sampling  P( Q | e )



- Gibbs Sampling  P( Q | e )