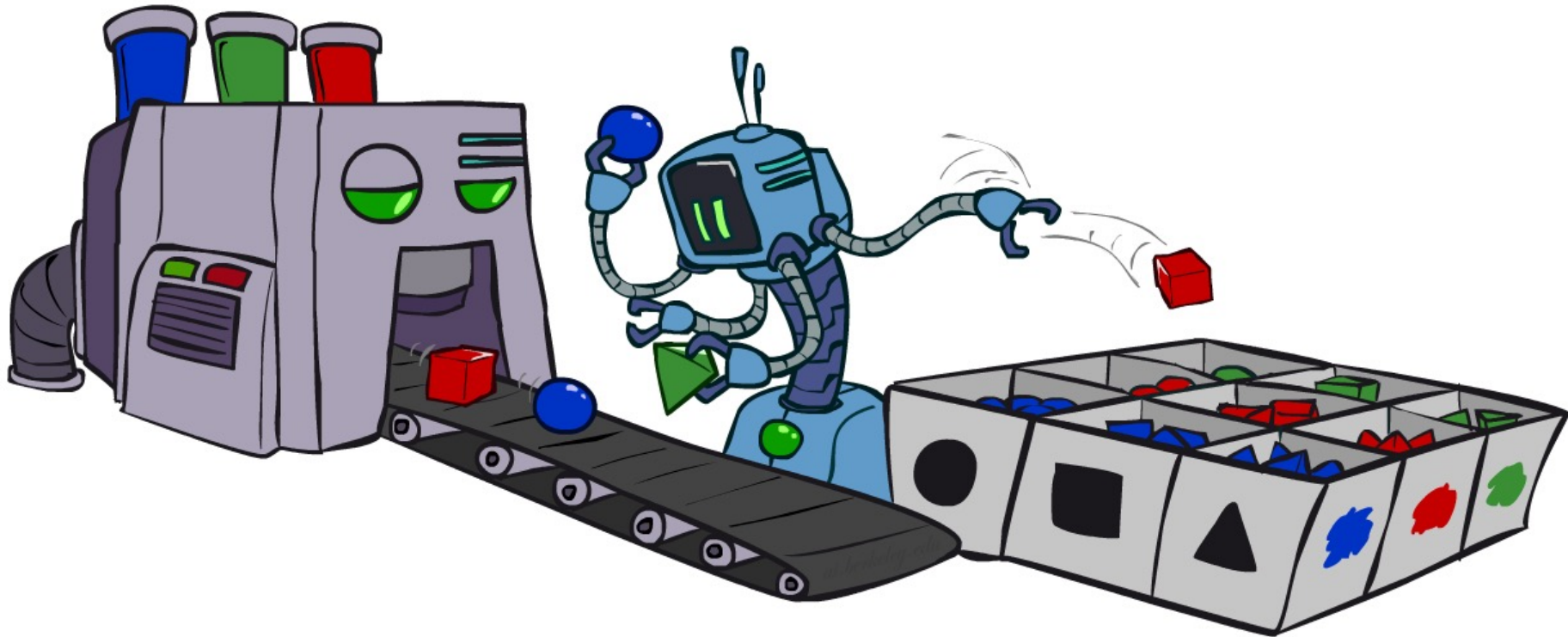


CS 3317: Artificial Intelligence

Bayes Nets: Sampling



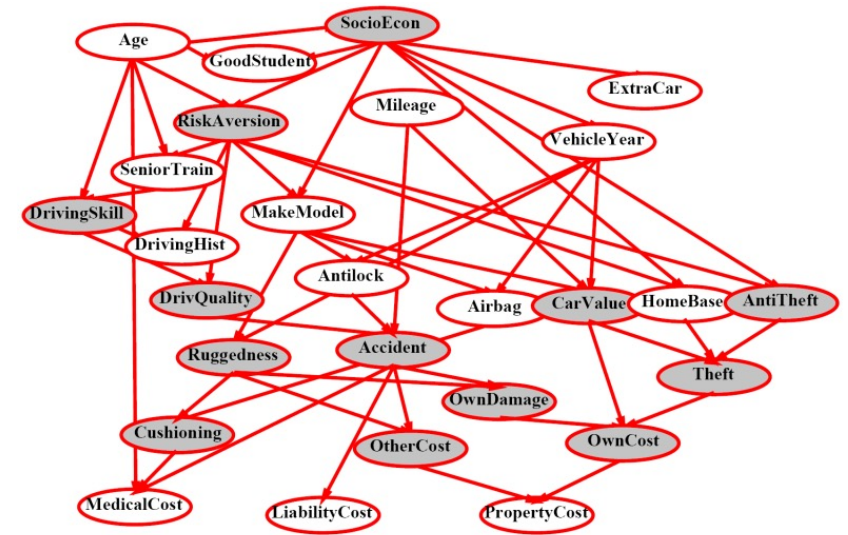
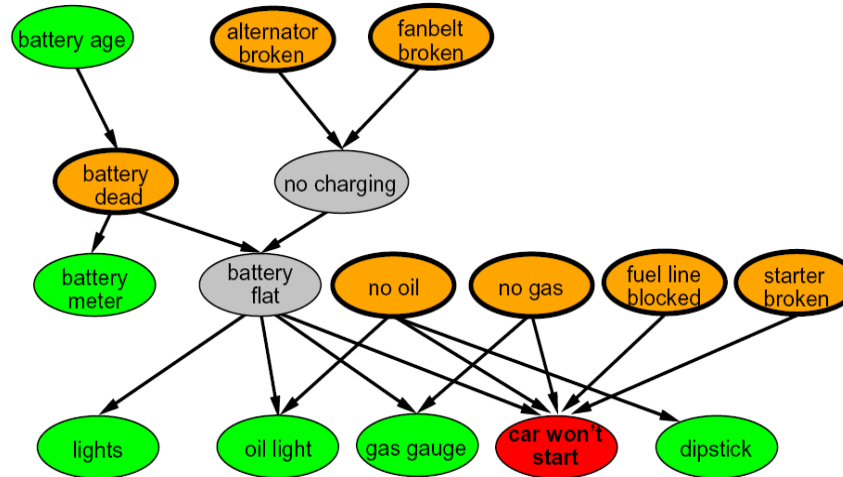
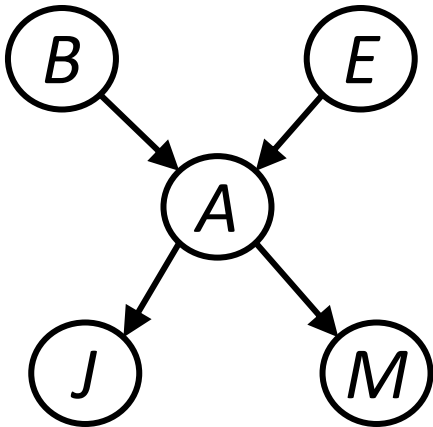
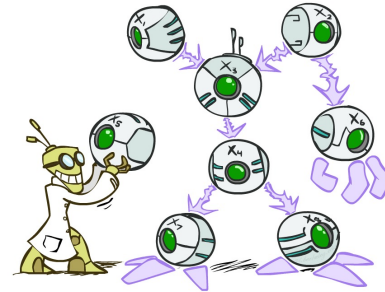
Instructor: Panpan Cai

[Slides adapted from UC Berkeley CS188]



Bayes Nets: Graphical Model

- A directed acyclic graph (DAG) consisting of random variables
 - A graphical model built from experience/data
- Examples:



Bayes Nets: Local CPTs

- A conditional probability table (CPT) for each node
 - A collection of distributions over X , one for each combination of parents' values

$$P(X|a_1 \dots a_n)$$

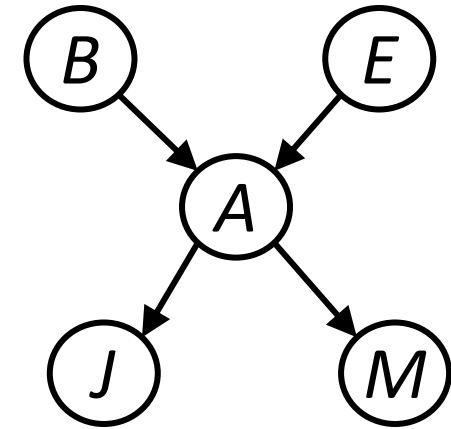
- Examples:

| B | P(B) |
|----|-------|
| +b | 0.001 |
| -b | 0.999 |

| E | P(E) |
|----|-------|
| +e | 0.002 |
| -e | 0.998 |

| A | J | P(J A) |
|----|----|--------|
| +a | +j | 0.9 |
| +a | -j | 0.1 |
| -a | +j | 0.05 |
| -a | -j | 0.95 |

| A | M | P(M A) |
|----|----|--------|
| +a | +m | 0.7 |
| +a | -m | 0.3 |
| -a | +m | 0.01 |
| -a | -m | 0.99 |



| B | E | A | P(A B,E) |
|----|----|----|----------|
| +b | +e | +a | 0.95 |
| +b | +e | -a | 0.05 |
| +b | -e | +a | 0.94 |
| +b | -e | -a | 0.06 |
| -b | +e | +a | 0.29 |
| -b | +e | -a | 0.71 |
| -b | -e | +a | 0.001 |
| -b | -e | -a | 0.999 |

Bayes Nets: Probabilities and Independences

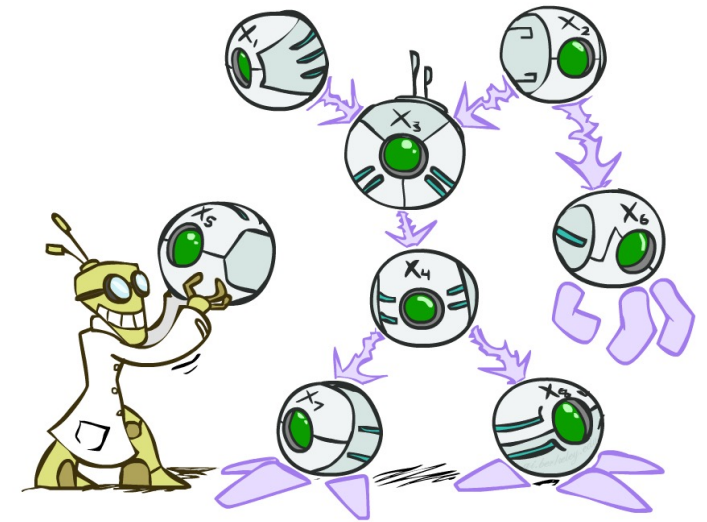
- Bayes nets implicitly encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

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

- The following conditional independence assumptions are made immediately, when given a BN graph:

$$P(x_i|x_1, \dots, x_{i-1}) = P(x_i|\text{parents}(X_i))$$

- X_i is independent of all non-parent variables coming before it, given its parents

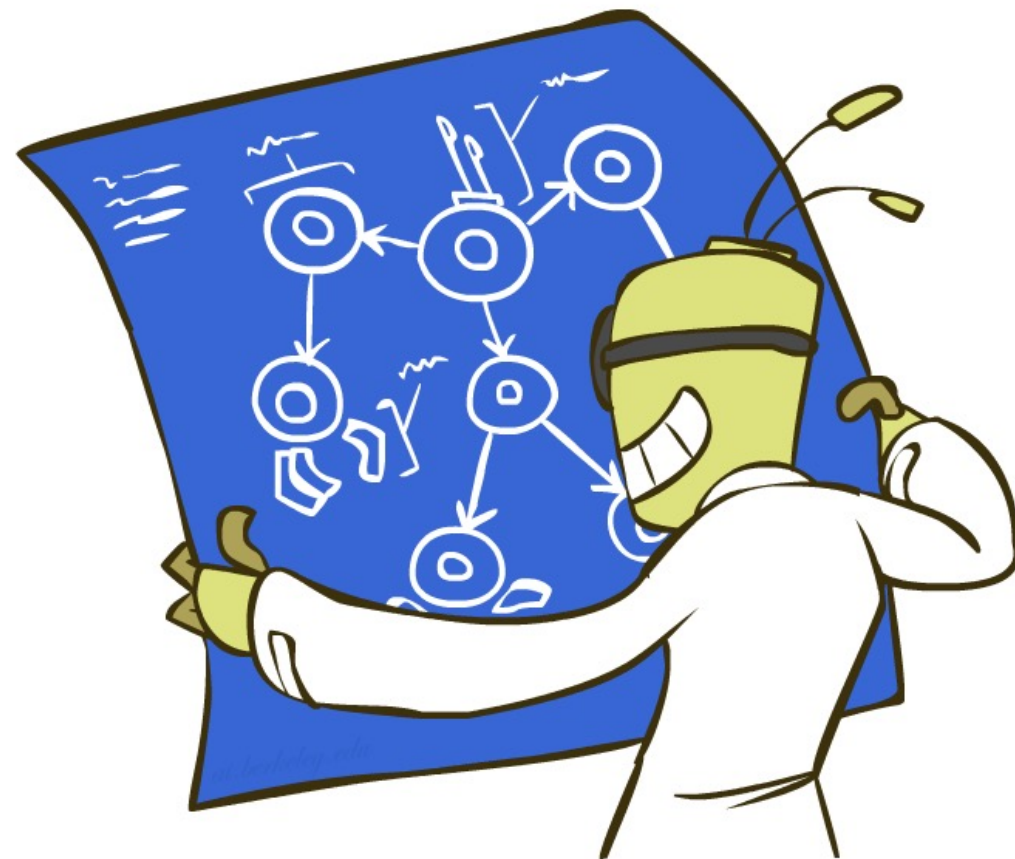


“Read” All Independences

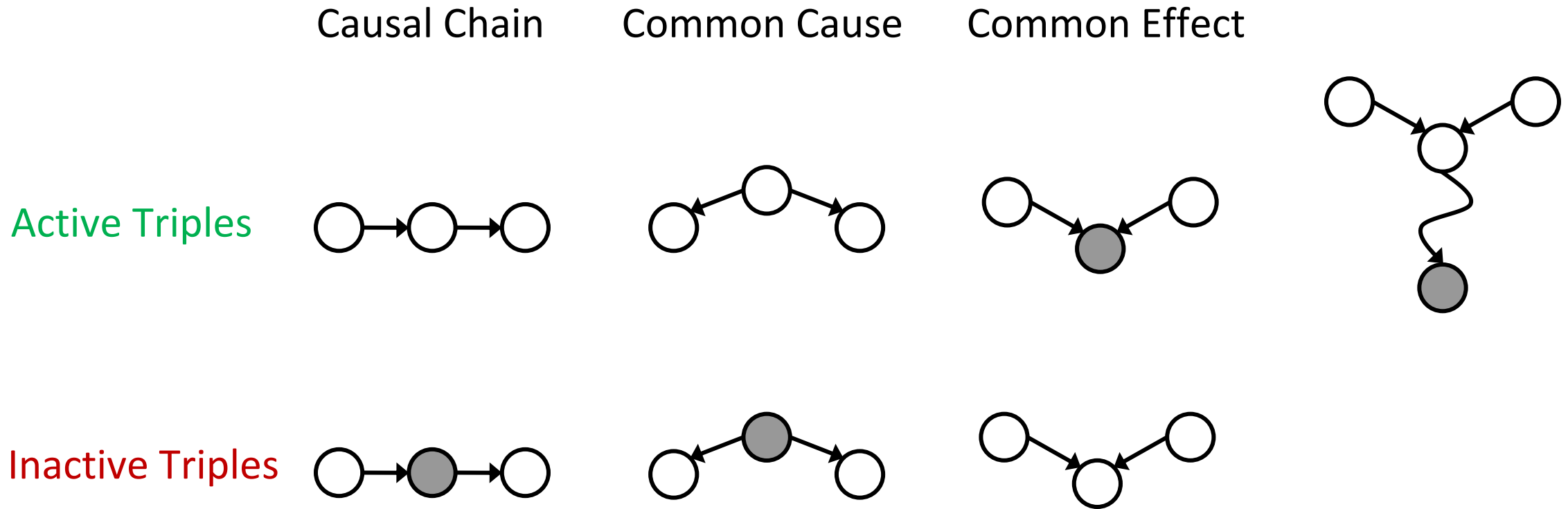
- Given a Bayes net structure, can run d-separation algorithm to read a *complete* list of conditional independences of the form:

$$X_i \perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$$

- This list determines the set of probability distributions that a BN graph can represent



D-Separation: Canonical Cases



D-Separation: Algorithm

- Query: $X_i \perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\} ?$

- Check all *undirected* paths between X_i and X_j

- If one or more path active, independence broken

$$X_i \not\perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$$

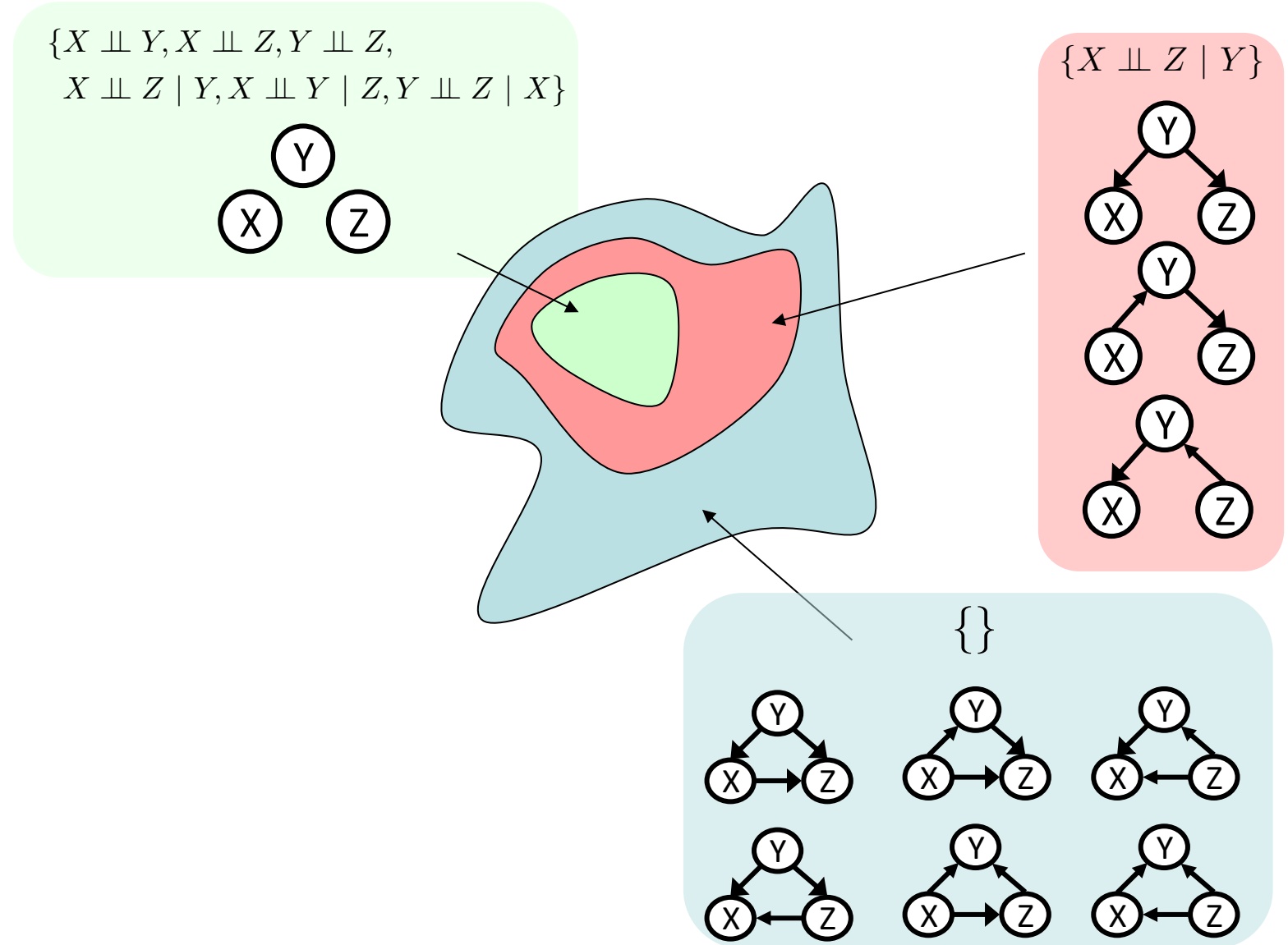
- Otherwise (i.e. if all paths are inactive), independence guaranteed

$$X_i \perp\!\!\!\perp X_j \mid \{X_{k_1}, \dots, X_{k_n}\}$$



Topology Limits Distributions

- Given some graph topology G , only a subset of joint distributions can be expressed
- Adding arcs increases the set of distributions, but at a cost of increased complexity
- Full conditioning can encode any distribution



Bayes Nets

- ✓ Representation
- ✓ Probabilistic Inference
 - Enumeration (exact, exponential complexity)
 - Variable elimination (exact, worst-case exponential complexity, often better)
 - Probabilistic inference is NP-complete
- ✓ Conditional Independences
- ✓ Sampling
 - Learning from data

Exact Inference: Inference by Enumeration

- Query: $P(Q|E_1 = e_1, \dots, E_k = e_k)$

- Start with initial factors:

- Local CPTs (but instantiated by evidence)

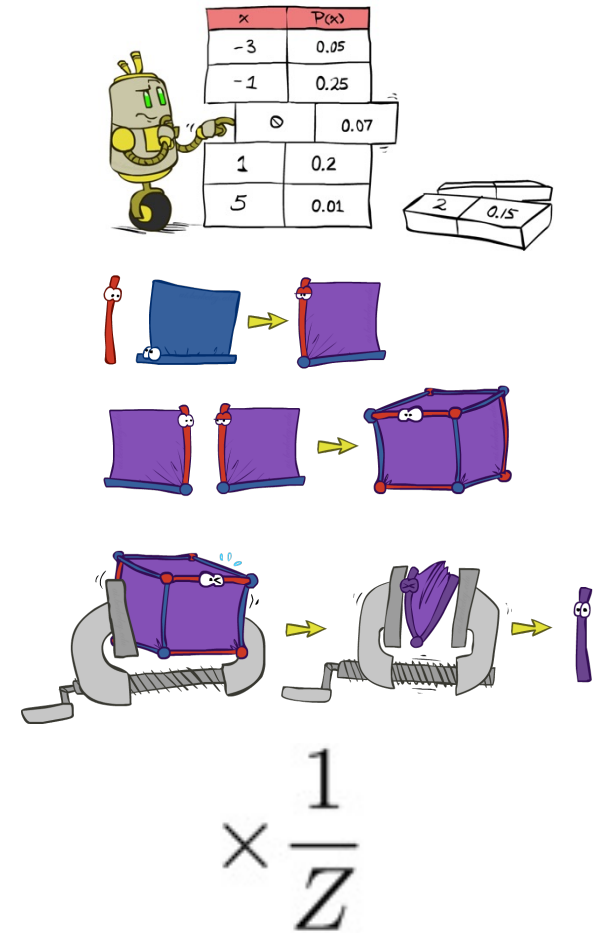
- Join all factors to get a complete joint

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

- Sum out all hidden variables


$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(Q, h_1 \dots h_r, e_1 \dots e_k)$$

- Normalize

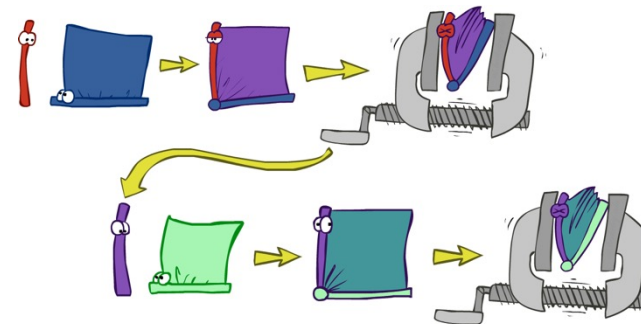


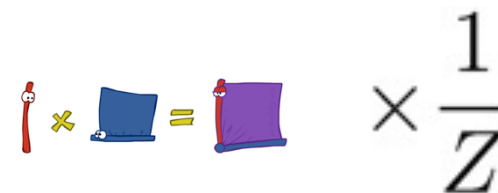
Exact Inference: Variable Elimination

- Query: $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- Start with initial factors:
 - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
 - Pick a hidden variable H
 - Join all factors mentioning H
 - Eliminate (sum out) H
- Join all remaining factors and normalize



| x | P(x) |
|----|------|
| -3 | 0.05 |
| -1 | 0.25 |
| 0 | 0.07 |
| 1 | 0.2 |
| 5 | 0.01 |




$$\text{red stick} \times \text{blue square} = \text{purple square} \times \frac{1}{Z}$$

Worst-Case Complexity: NP Hard!

- CSP: assign variables so that the sentence is **true**

Clause need be **true**

$$(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_3 \vee \neg x_4) \wedge (x_2 \vee \neg x_2 \vee x_4) \wedge (\neg x_3 \vee \neg x_4 \vee \neg x_5) \wedge (x_2 \vee x_5 \vee x_7) \wedge (x_4 \vee x_5 \vee x_6) \wedge (\neg x_5 \vee x_6 \vee \neg x_7) \wedge (\neg x_5 \vee \neg x_6 \vee x_7)$$

$$P(X_i = 0) = P(X_i = 1) = 0.5$$

$$Y_1 = X_1 \vee X_2 \vee \neg X_3$$

...

$$Y_8 = \neg X_5 \vee X_6 \vee X_7$$

$$Y_{1,2} = Y_1 \wedge Y_2$$

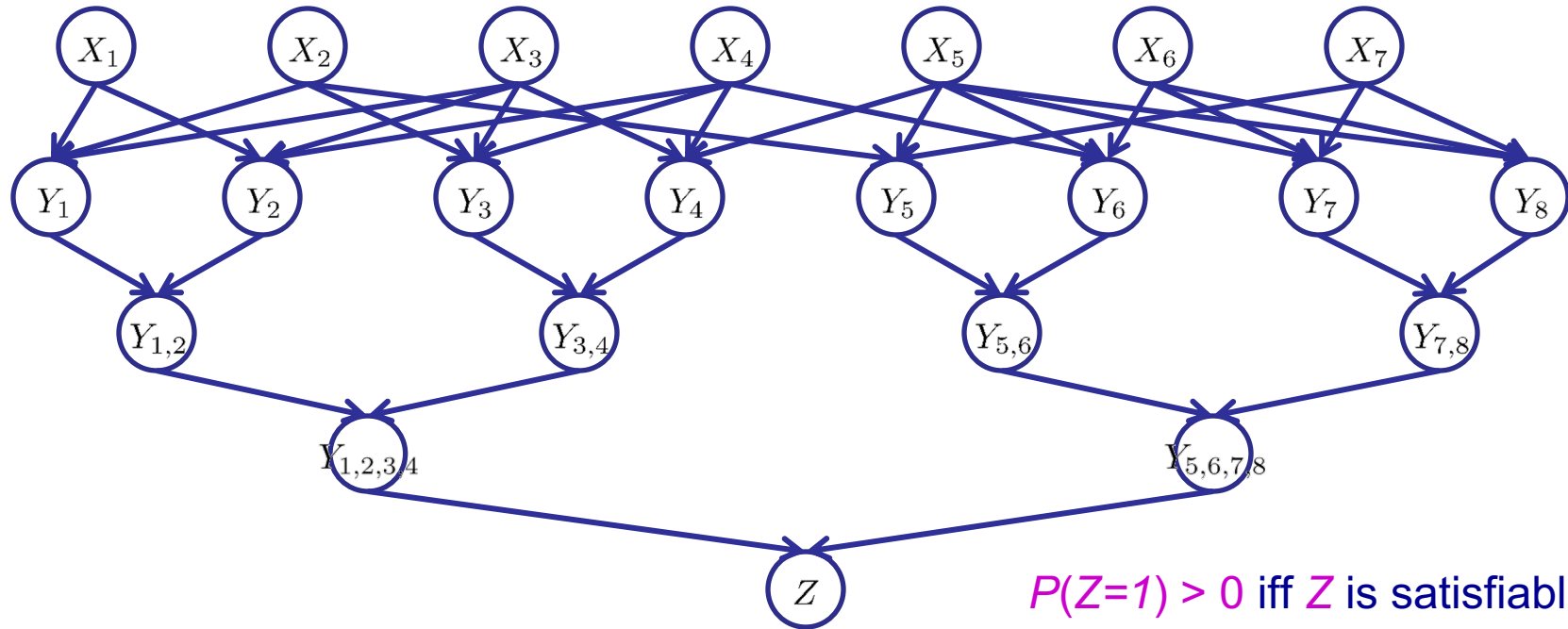
...

$$Y_{7,8} = Y_7 \wedge Y_8$$

$$Y_{1,2,3,4} = Y_{1,2} \wedge Y_{3,4}$$

$$Y_{5,6,7,8} = Y_{5,6} \wedge Y_{7,8}$$

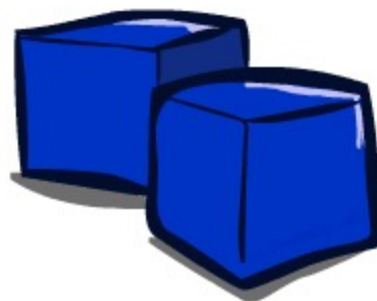
$$Z = Y_{1,2,3,4} \wedge Y_{5,6,7,8}$$



$P(Z=1) > 0$ iff Z is satisfiable

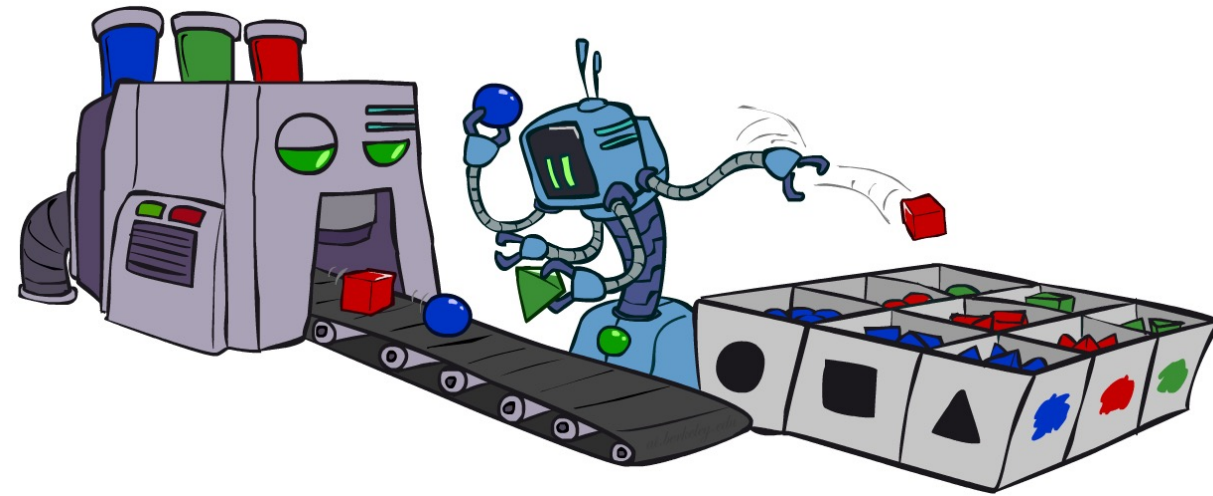
- If we can answer whether $P(z)$ equal to zero or not, we answered whether the 3-SAT problem has a solution (which is NP hard!)
- Hence probabilistic inference in Bayes nets is NP-hard.
 - No known efficient (exact) solution *in general*.

Approximate Inference: Sampling



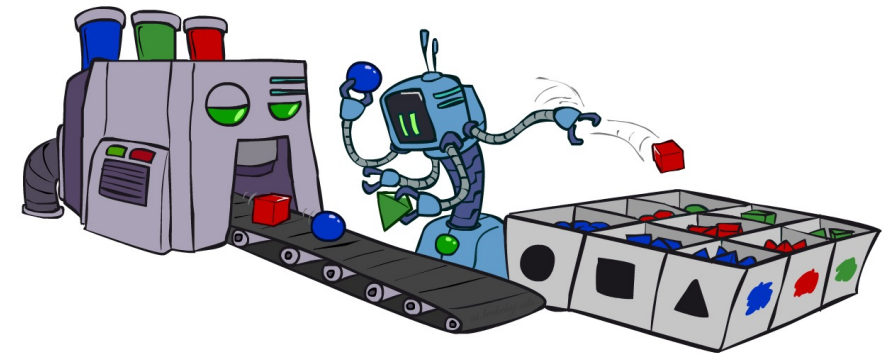
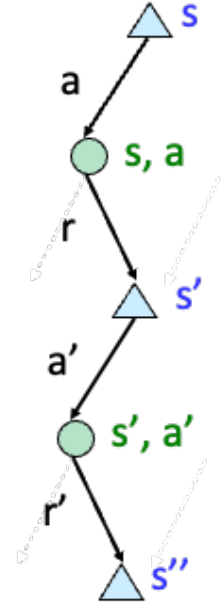
Sampling

- Sampling is just repeated random simulation
 - X is a random variable with an underlying distribution $P(X)$
 - Treat the model $P(X)$ as a simulator, where each simulation generates a sample x_i
 - In the limit, the *frequency* of seeing a particular value x converges to the *probability* of x
- The simulator can also host:
 - Joint distributions
 - Conditional distributions
 - Bayes Nets (with evidences or not)
 - Or any *black-box* random process...



Sampling

- Why sample?
 - *Learning*: estimate a distribution that you don't know
 - *Model-based RL*: *sample* transitions, learn models
 - *Q-learning*: *sample* transitions, learn *Q*-functions
 - *Speedup inference*: estimating using samples is faster than computing the exact answer (e.g. with variable elimination)
 - *This lecture!*



How to Sample?

- Sampling from given distribution

- Step 1: Get sample u from uniform distribution over $[0, 1)$
 - E.g. `random()` in python
- Step 2: Convert u into an outcome
 - Associate each outcome with a *sub-interval* of $[0,1)$ with *size* equal to *probability* of the outcome

- Example

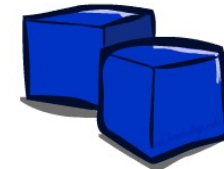
| C | P(C) |
|-------|------|
| red | 0.6 |
| green | 0.1 |
| blue | 0.3 |

$$0 \leq u < 0.6, \rightarrow C = \text{red}$$

$$0.6 \leq u < 0.7, \rightarrow C = \text{green}$$

$$0.7 \leq u < 1, \rightarrow C = \text{blue}$$

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



Sampling in Bayes Nets

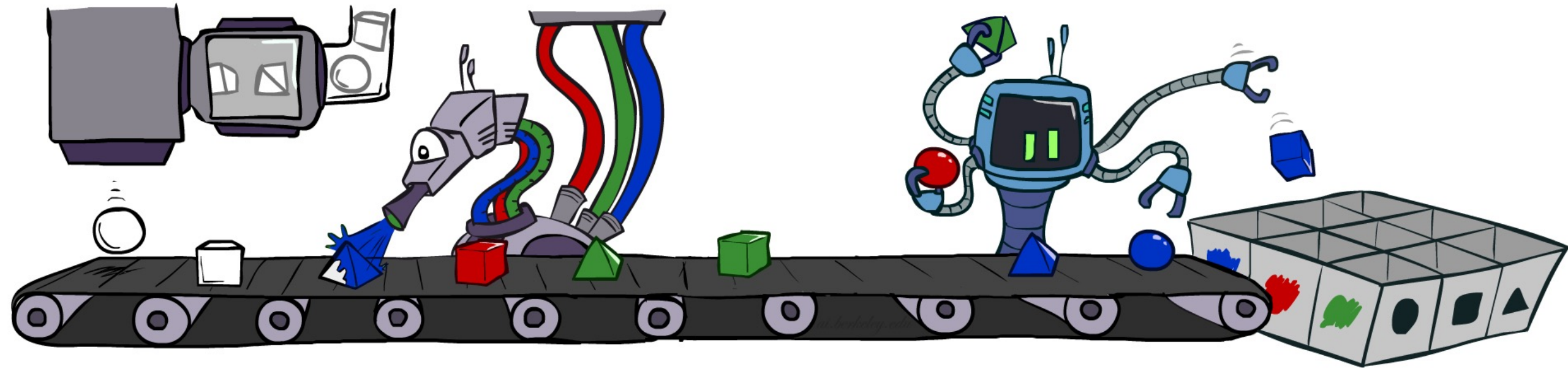
- Prior Sampling
- Rejection Sampling
- Likelihood Weighting
- Gibbs Sampling

Prior Sampling

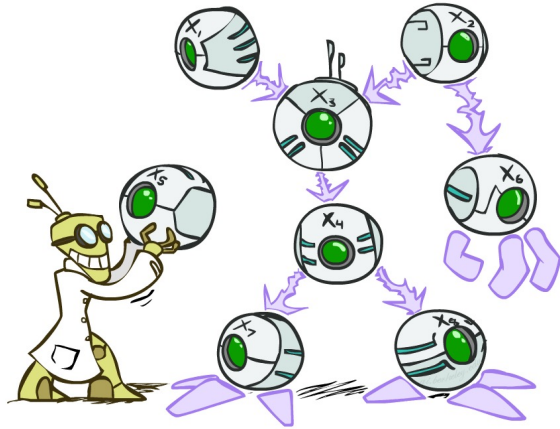
"Shaping machine"

"Coloring machine"

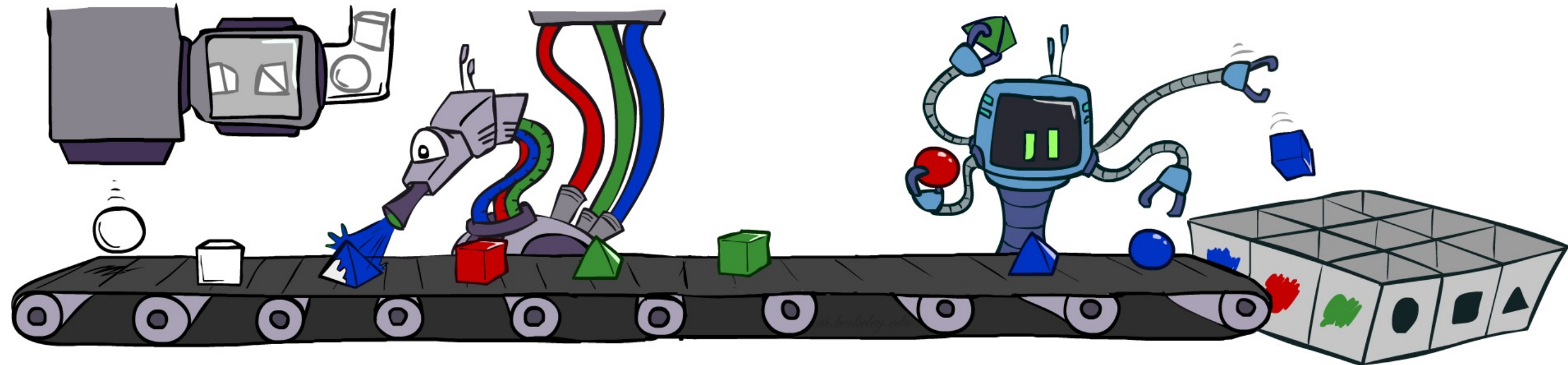
"Sorting machine"



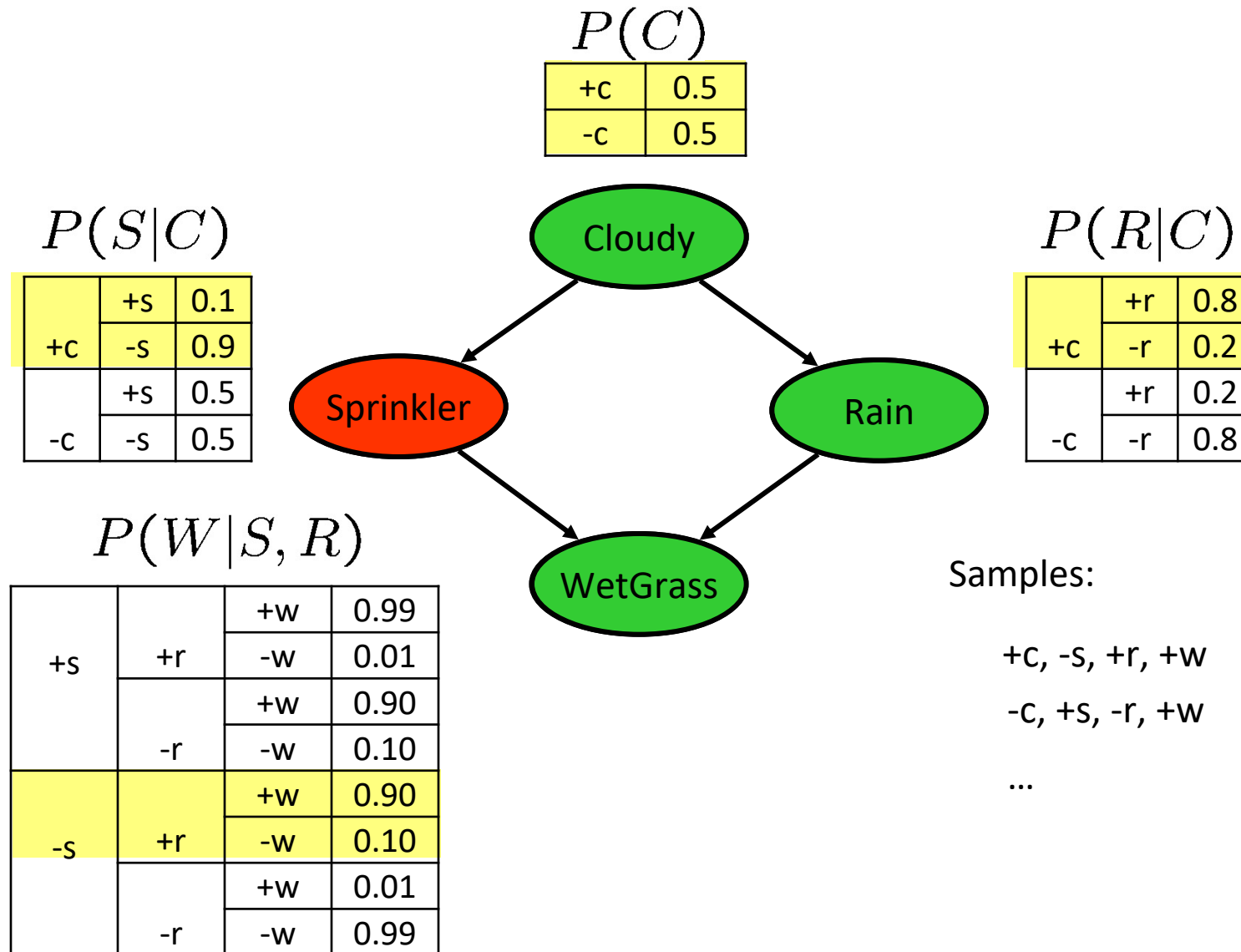
Prior Sampling



- For $i = 1, 2, \dots, n$ in topological order
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
- Return (x_1, x_2, \dots, x_n)



Example: Prior Sampling



Prior Sampling

- This process generates samples with probability:

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


*Probability of independent
events happening together*

*Product of probability
of each event*

BN's joint probability

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

Then $\lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) = \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N$

$$= S_{PS}(x_1, \dots, x_n)$$
$$= P(x_1 \dots x_n)$$


- I.e., the sampling procedure is **consistent**

Example: Prior Sampling

- 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

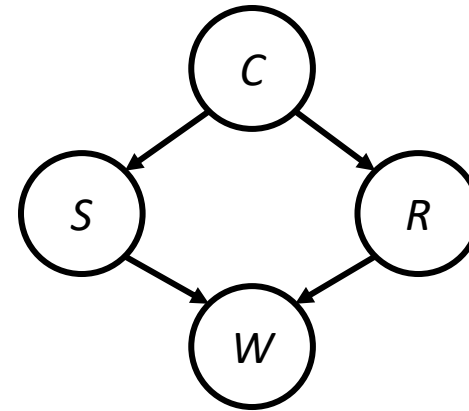
- If we want to know $P(W)$

- We have counts $\langle +w:4, -w:1 \rangle$
- Normalize to get $P(W) = \langle +w:0.8, -w:0.2 \rangle$
- This will get closer to the true distribution with more samples

- Can estimate anything else, too

- $P(C \mid +w)$? $P(C \mid +r, +w)$?
- What about $P(C \mid -r, -w)$?

- Can also use this to estimate expected value of a function $E[f(X)]$ - *Monte Carlo Estimation*



Connect: Inference by Enumeration

- Query: $P(Q|E_1 = e_1, \dots, E_k = e_k)$

- Start with initial factors:
 - Local CPTs (but instantiated by evidence)



2. Select only relevant samples

- Join all factors to get a complete joint

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



1. Sample complete values

- Sum out all hidden variables

$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(Q, h_1 \dots h_r, e_1 \dots e_k)$$



3. Counting (sum up weights of samples)

- Normalize



4. Normalize

Problem of Prior Sampling

- 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

Hosting all samples here!

- If we want to know $P(W)$

- We have counts $\langle +w:4, -w:1 \rangle$
- Normalize to get $P(W) = \langle +w:0.8, -w:0.2 \rangle$
- This will get closer to the true distribution with more samples

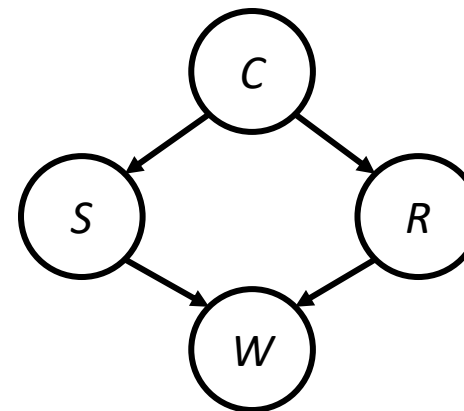
- Can estimate anything else, too

- $P(C \mid +w)$? $P(C \mid +r, +w)$?

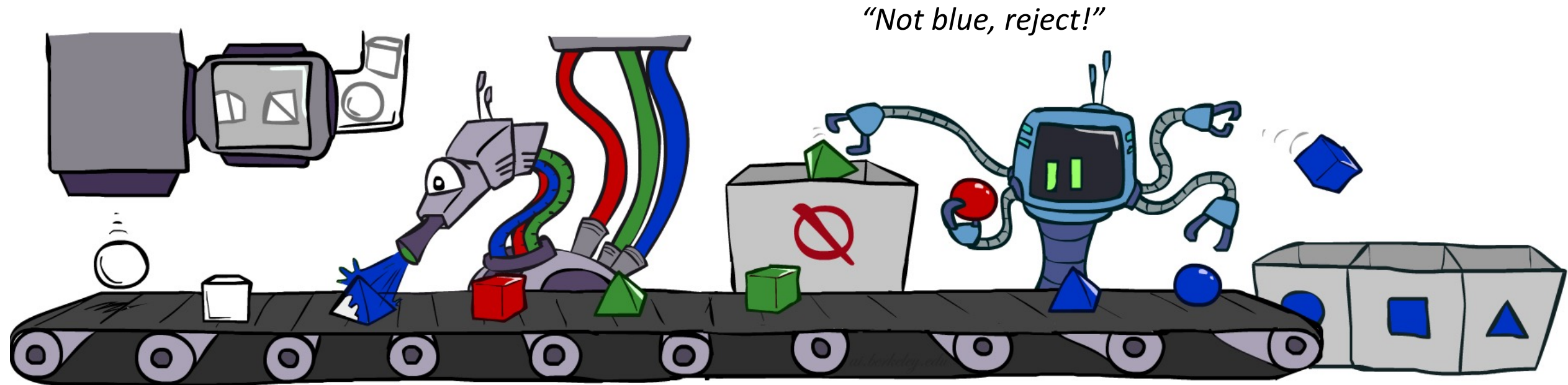
- What about $P(C \mid -r, -w)$?

- *Many samples are not consistent with evidence*

- *What we need are just numbers, e.g., $N_{ps}(+c \mid +r, +w)$*

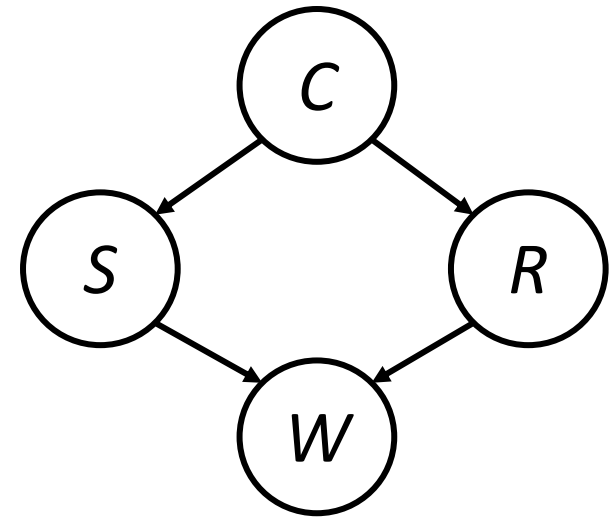


Rejection Sampling



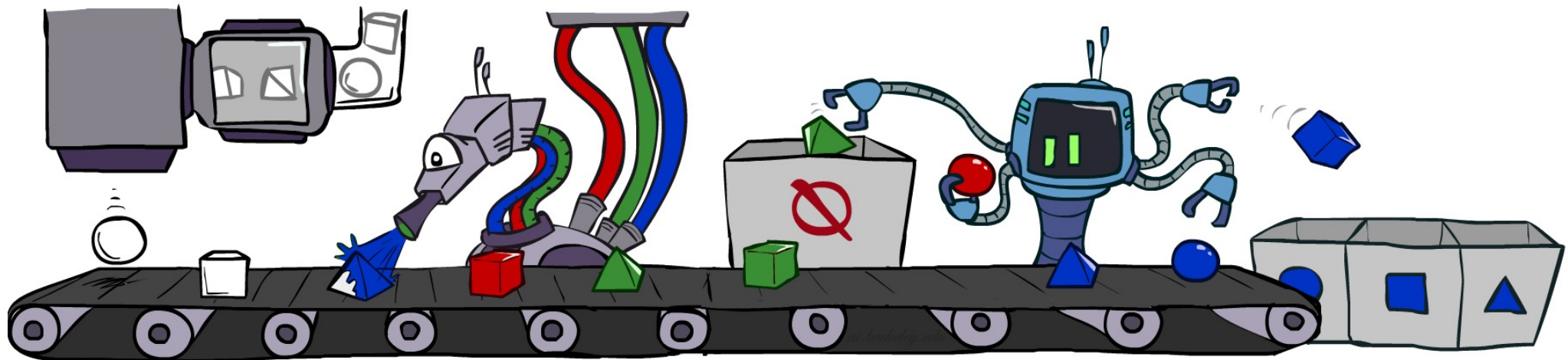
Rejection Sampling

- Let's say we want $P(C)$
 - Nothing to reject, just track counts of c as we go (no need to store samples)
- Let's say we want $P(C \mid +s)$
 - Count c , but reject samples which don't have $S=+s$
 - This is called *rejection sampling*
 - Benefit: we can filter out samples early!
 - Property: also *consistent* for conditional probabilities (i.e., correct in the limit)



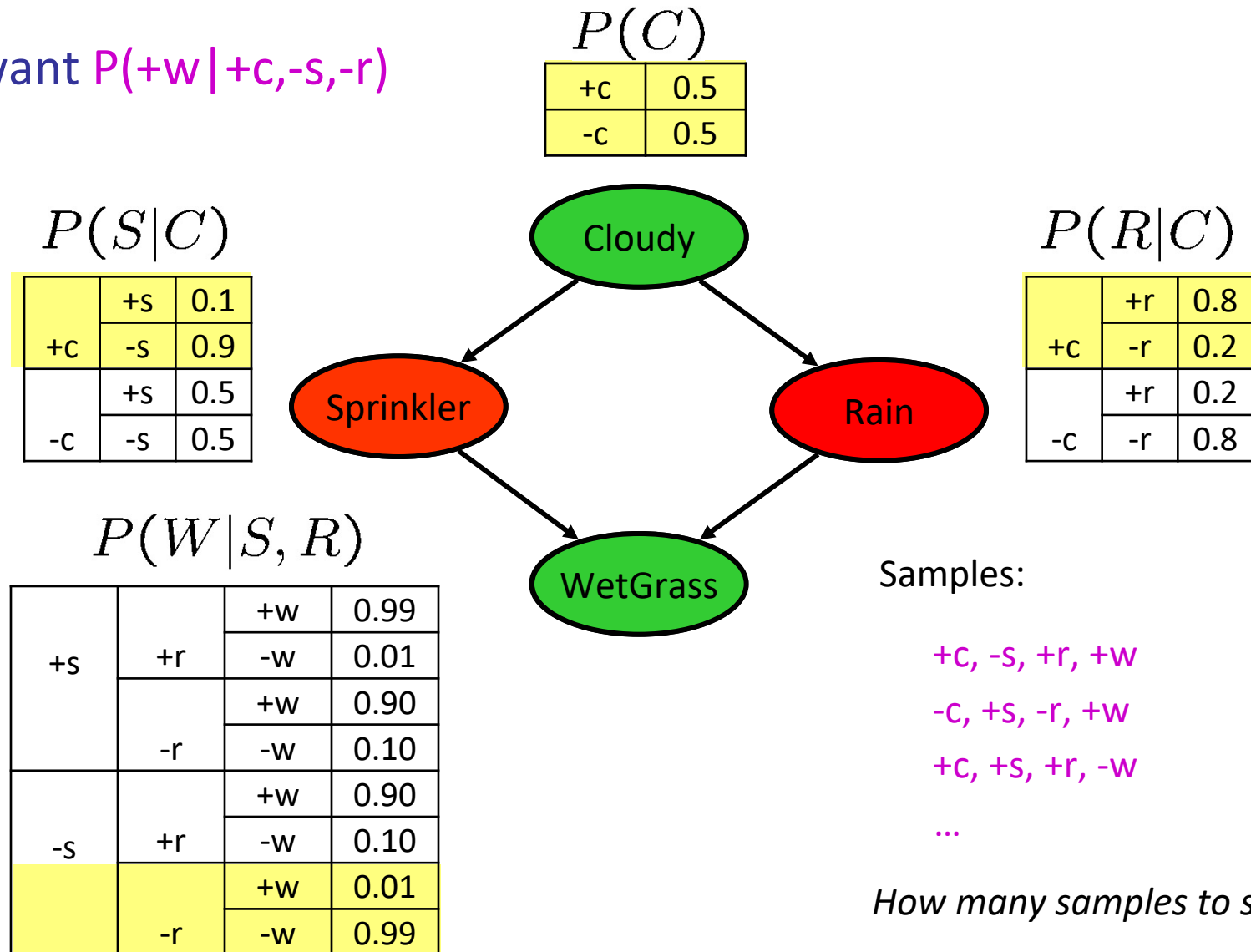
Rejection Sampling

- Input: evidence instantiation
- For $i = 1, 2, \dots, n$ in topological order
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
 - If x_i not consistent with evidence
 - Return (*reject*: no sample is generated in this cycle)
- Return (x_1, x_2, \dots, x_n)



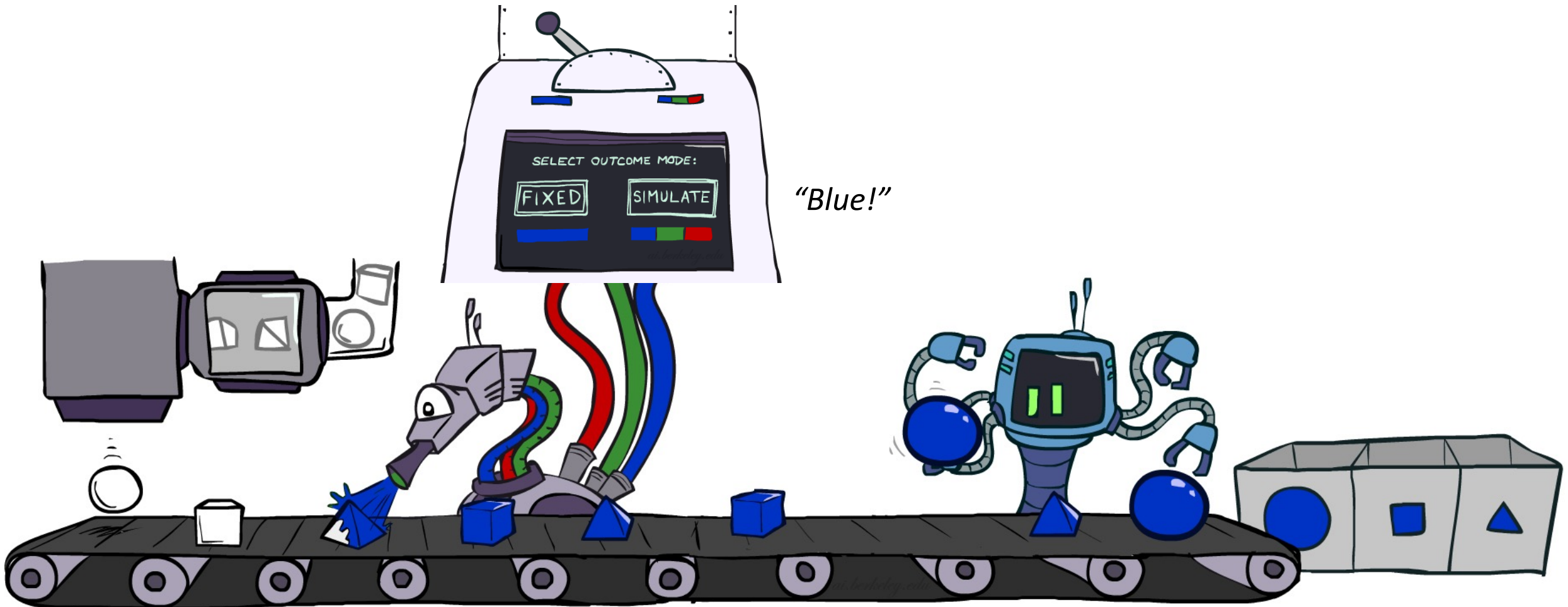
Problem of Rejection Sampling

- Let's say we want $P(+w | +c, -s, -r)$



Likelihood Weighting

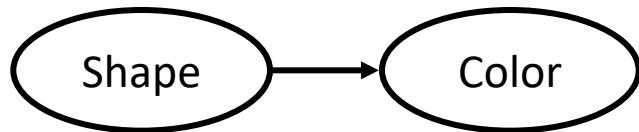
"Hacked machine"



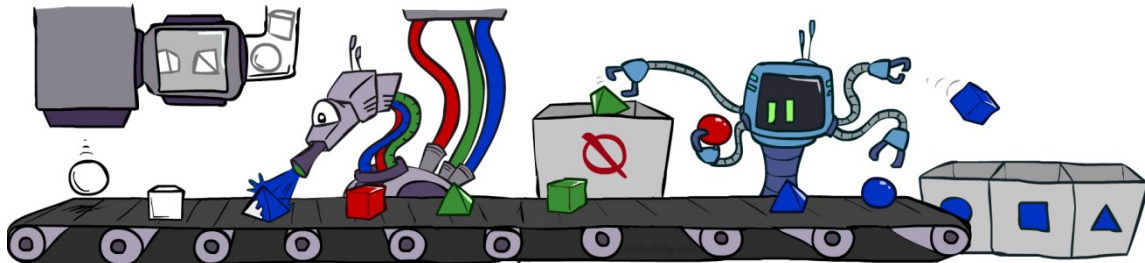
"Blue!"

Likelihood Weighting

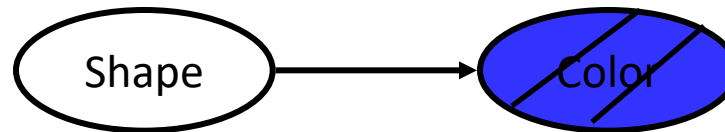
- Problem with rejection sampling:
 - If evidence is unlikely, rejects lots of samples
 - Consider $P(\text{Shape} \mid \text{blue})$



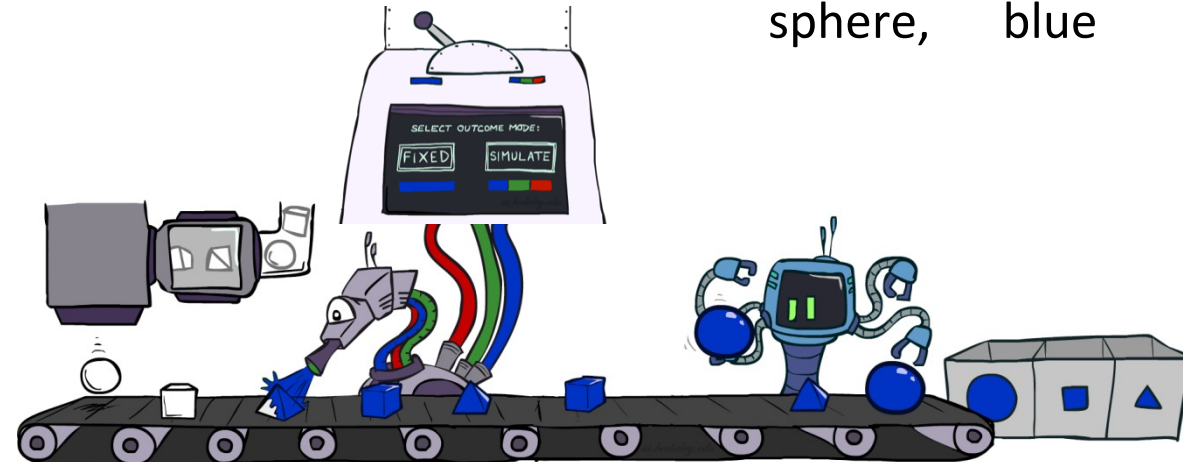
pyramid, green
pyramid, red
sphere, blue
cube, red
sphere, green



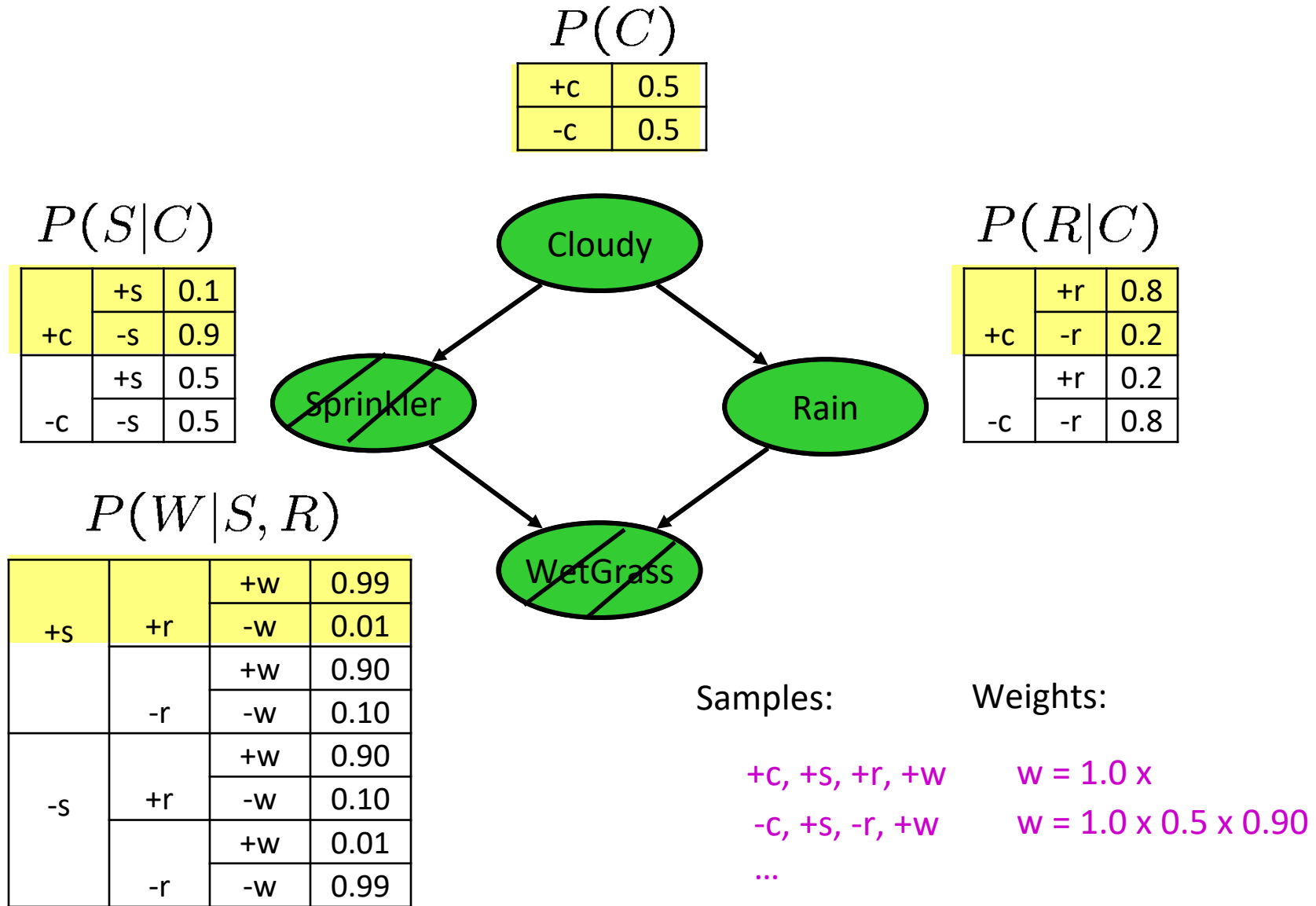
- Idea: *fix* evidence variables and sample the rest
 - Problem: sample distribution not consistent!
 - Solution: *weight* by probability of evidence given parents



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



Likelihood Weighting



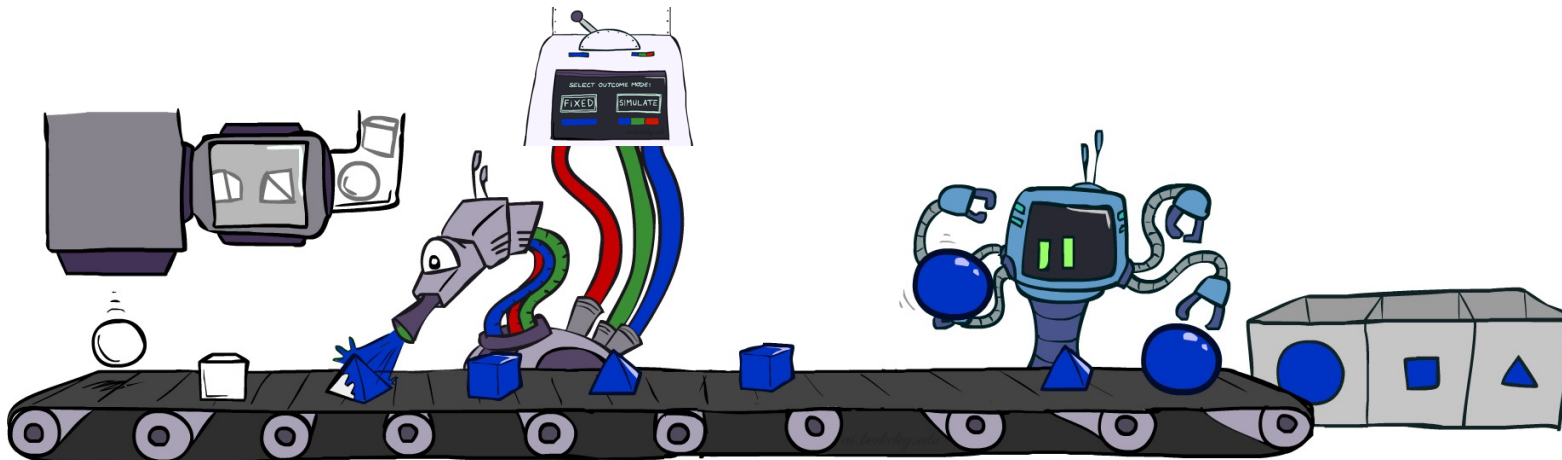
Likelihood Weighting

- Input: evidence instantiation
- $w = 1.0$
- for $i = 1, 2, \dots, n$ in topological order
 - if X_i is an evidence variable
 - $X_i = \text{observation } x_i$
 - 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 $(x_1, x_2, \dots, x_n), w$

Counting



Summing up weights



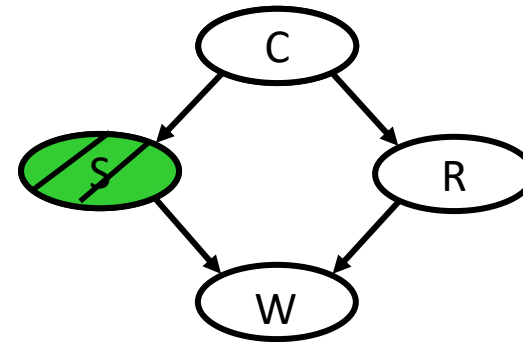
Likelihood Weighting

- Sampling distribution if \mathbf{z} sampled and evidence \mathbf{e} fixed

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

- Now, correct with weights

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



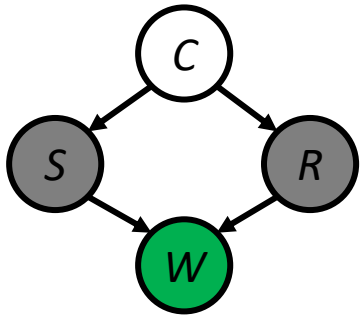
- Combining together, weighted sampling distribution becomes *consistent*

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

Likelihood Weighting

- Likelihood weighting is helpful

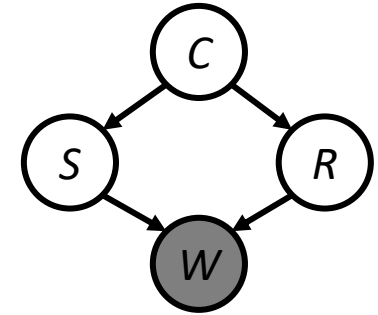
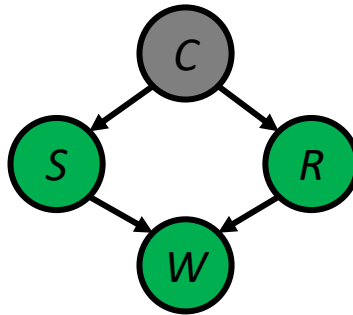
- We have taken evidence into account as we generate the sample
 - *E.g. here, W 's value will be picked based on fixed evidence values of S, R*
- More of our samples reflect the state of the world suggested by the evidence



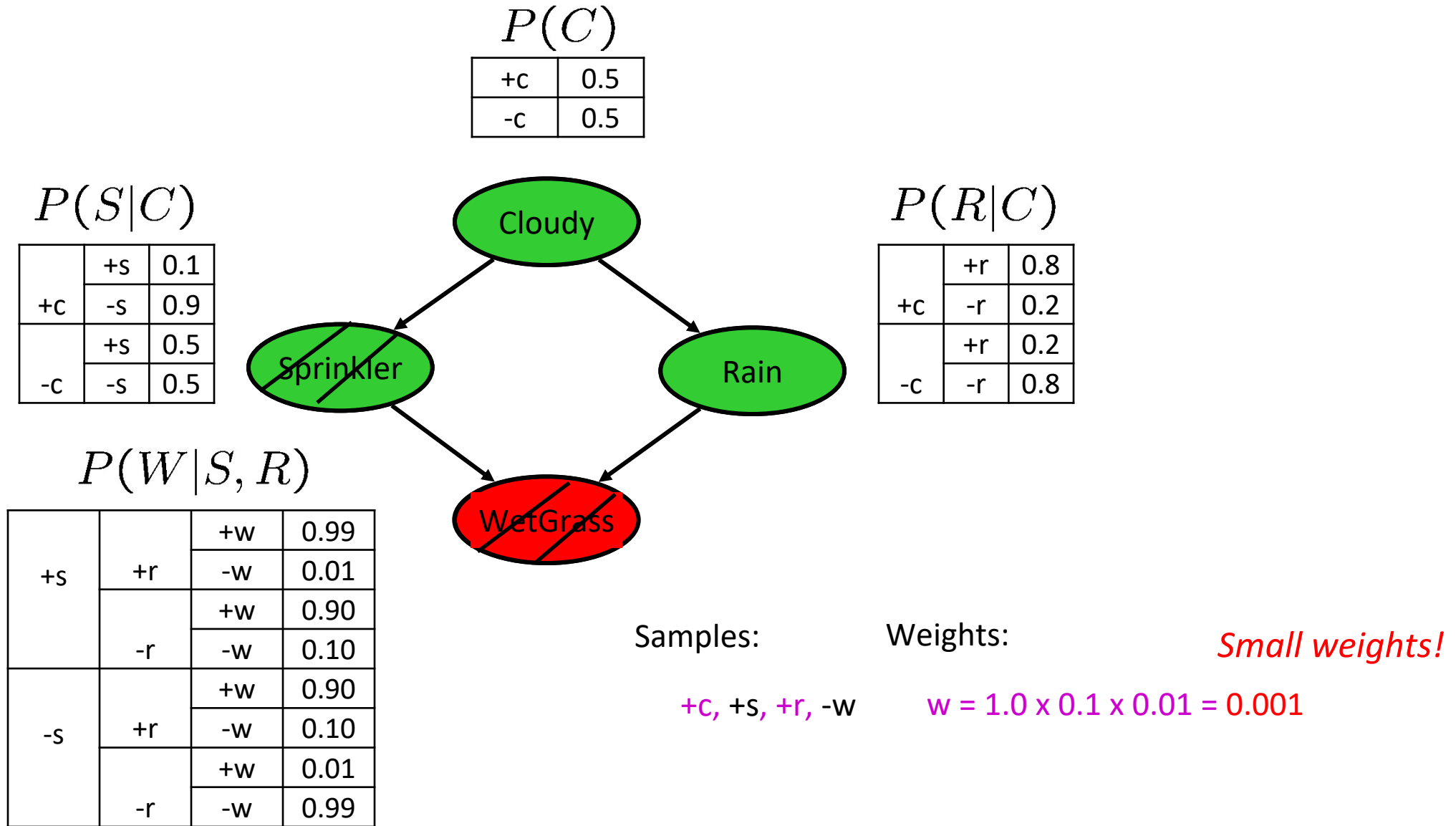
- Likelihood weighting doesn't solve all our problems

- Fixation of evidence influences the choice of *downstream* variables, but not *upstream* ones
 - W is more likely to match the evidence c
 - C isn't more likely to match the evidence w

- We would like to consider the influence of fixed evidence when sampling all variables (Gibbs sampling!)



Problem of Likelihood Weighting



Gibbs Sampling

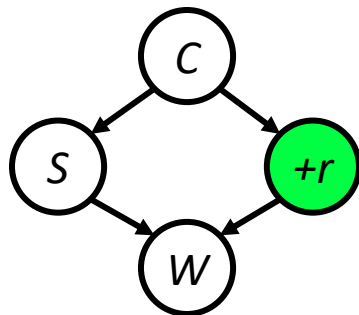
"Iterative update!"



Gibbs Sampling Example: $P(S \mid +r)$

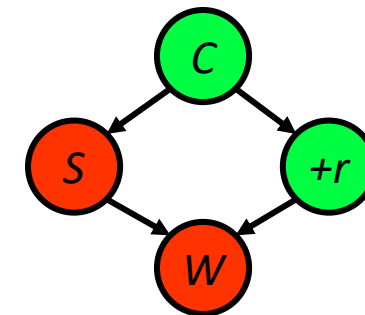
- Step 1: Fix evidence

- $R = +r$



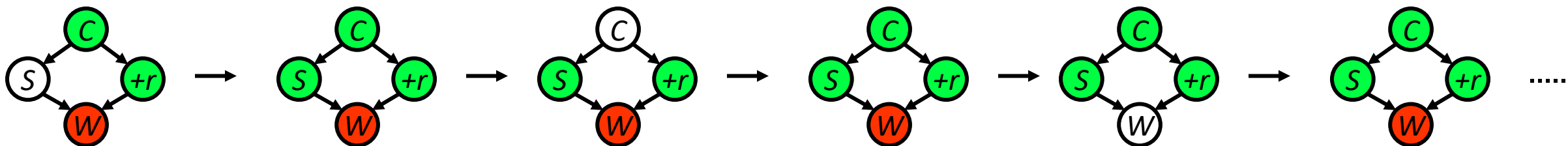
- Step 2: Initialize other variables

- Randomly



- Steps 3: Repeat

- Choose a non-evidence variable X
- Resample X from $P(X \mid \text{all other variables})^*$



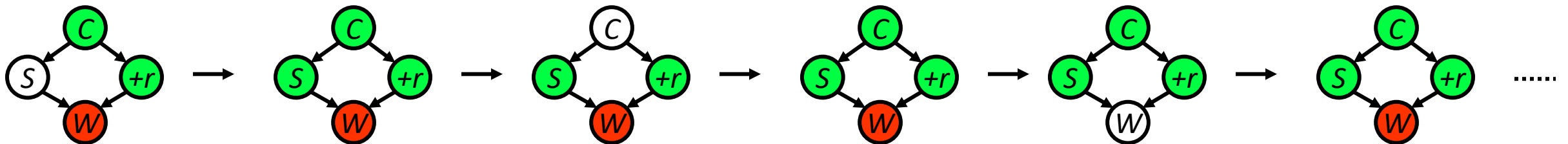
Sample from $P(S \mid +c, -w, +r)$

Sample from $P(C \mid +s, -w, +r)$

Sample from $P(W \mid +s, +c, +r)$

Gibbs Sampling

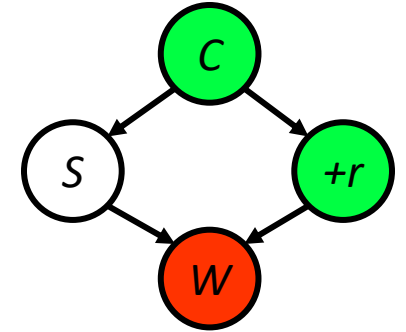
- **Procedure:** keep track of a *full instantiation* x_1, x_2, \dots, x_n .
 - Start with an *arbitrary* instantiation consistent with the evidence
 - Resample *one* (non-evidence) variable at a time, conditioned on *all the rest*
 - Keep repeating this for a long time
- **Property:** in the limit (repeated unlimited times), the resulting set of samples reflect the correct distribution (i.e. $P(Q|e)$).



Resampling of One Variable

- Sample from $P(S \mid +c, +r, -w)$

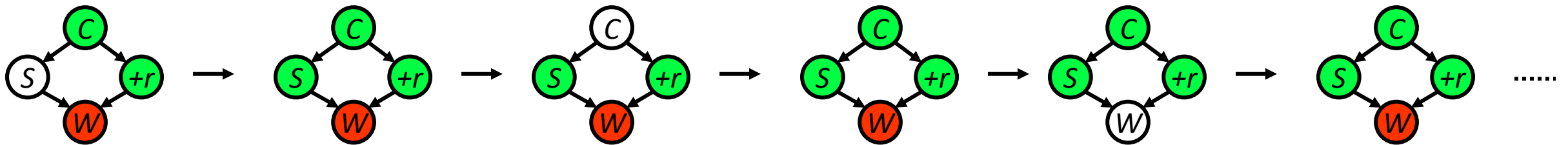
$$\begin{aligned} P(S \mid +c, +r, -w) &= \frac{P(S, +c, +r, -w)}{P(+c, +r, -w)} \\ &= \frac{P(S, +c, +r, -w)}{\sum_s P(s, +c, +r, -w)} \\ &= \frac{P(+c)P(S \mid +c)P(+r \mid +c)P(-w \mid S, +r)}{\sum_s P(+c)P(s \mid +c)P(+r \mid +c)P(-w \mid s, +r)} \\ &= \frac{P(+c)P(S \mid +c)P(+r \mid +c)P(-w \mid S, +r)}{P(+c)P(+r \mid +c) \sum_s P(s \mid +c)P(-w \mid s, +r)} \\ &= \frac{P(S \mid +c)P(-w \mid S, +r)}{\sum_s P(s \mid +c)P(-w \mid s, +r)} \end{aligned}$$



- Many things cancel out – only CPTs with **S** remain!
- More generally: only CPTs with *the resampled variable* need to be considered
 - Join them together and normalize

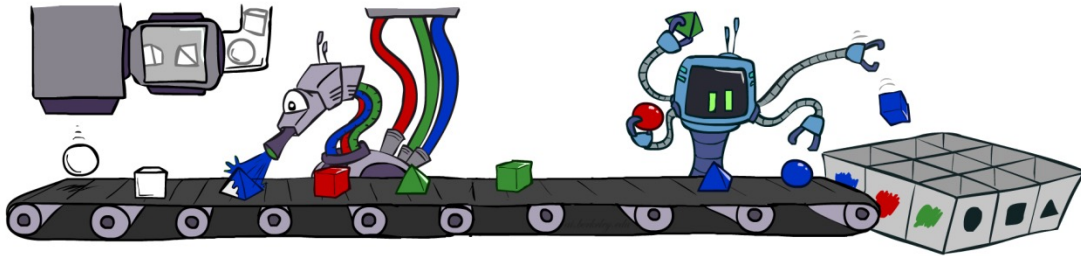
Gibbs Sampling

- **Rationale:** pick both *upstream* and *downstream* variables condition on evidence
 - All samples match the evidence, and are equally “effective” (with same weights)
 - Avoided small weights in likelihood weighting!

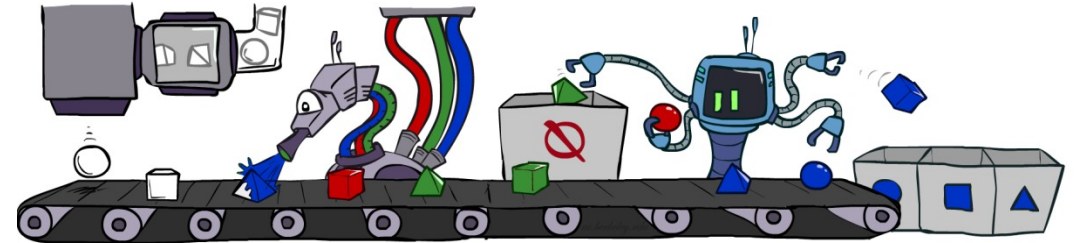


Bayes Net Sampling Summary

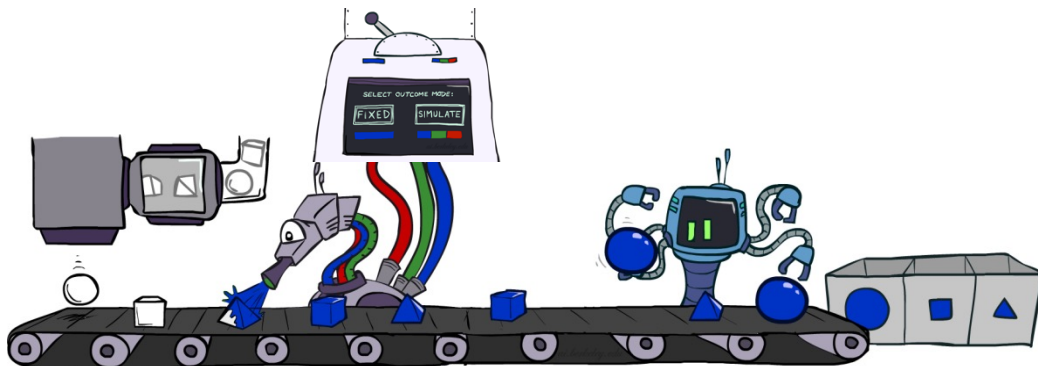
- Prior Sampling $P(Q)$



- Rejection Sampling $P(Q | e)$



- Likelihood Weighting $P(Q | e)$



- Gibbs Sampling $P(Q | e)$

