# Neural Particle Smoothing for Sampling from Conditional Sequence Models
## Chu-Cheng Lin and Jason Eisner — JOHNS HOPKINS UNIVERSITY

## Quiz question: CRF is to dynamic programming as RNN is to …?

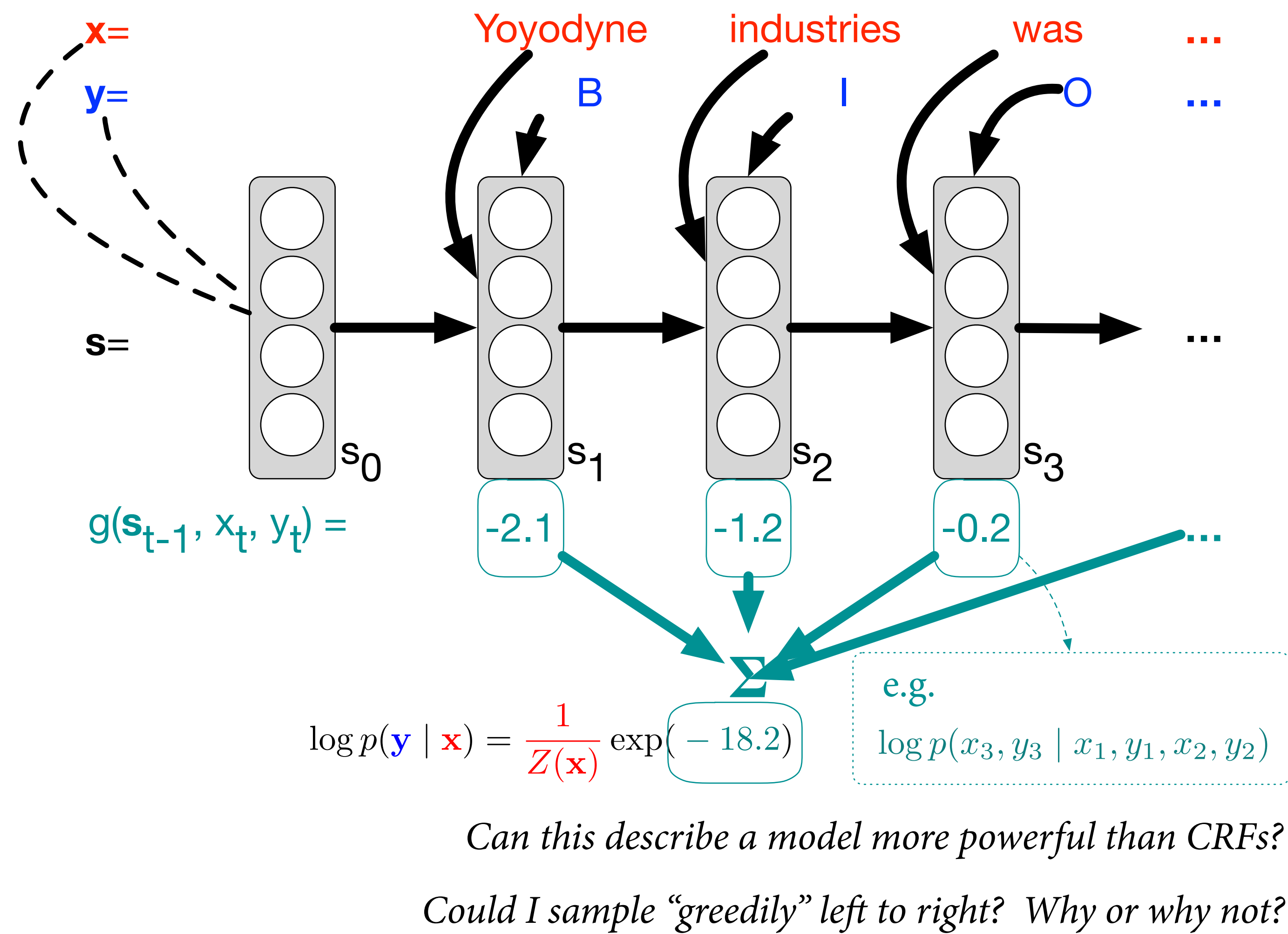Hi, I built a fancy globally normalized neural model that is more powerful than CRFs! 😊

But oops, dynamic programming doesn't work anymore 😟

— *How will I compute my gradient for training?*

— *How will I figure out what my model predicts (conditioned on evidence)?*

— *How will I combine my model with other probability distributions?*

I guess I can sample from $p(\mathbf{y} \mid \mathbf{x})$, but how? My model specifies $\tilde{p}(\mathbf{x}, \mathbf{y})$

I know that $p(\mathbf{y} \mid \mathbf{x}) = \frac{\tilde{p}(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}'} \tilde{p}(\mathbf{x}, \mathbf{y}')}$, but I don't really feel like summing over

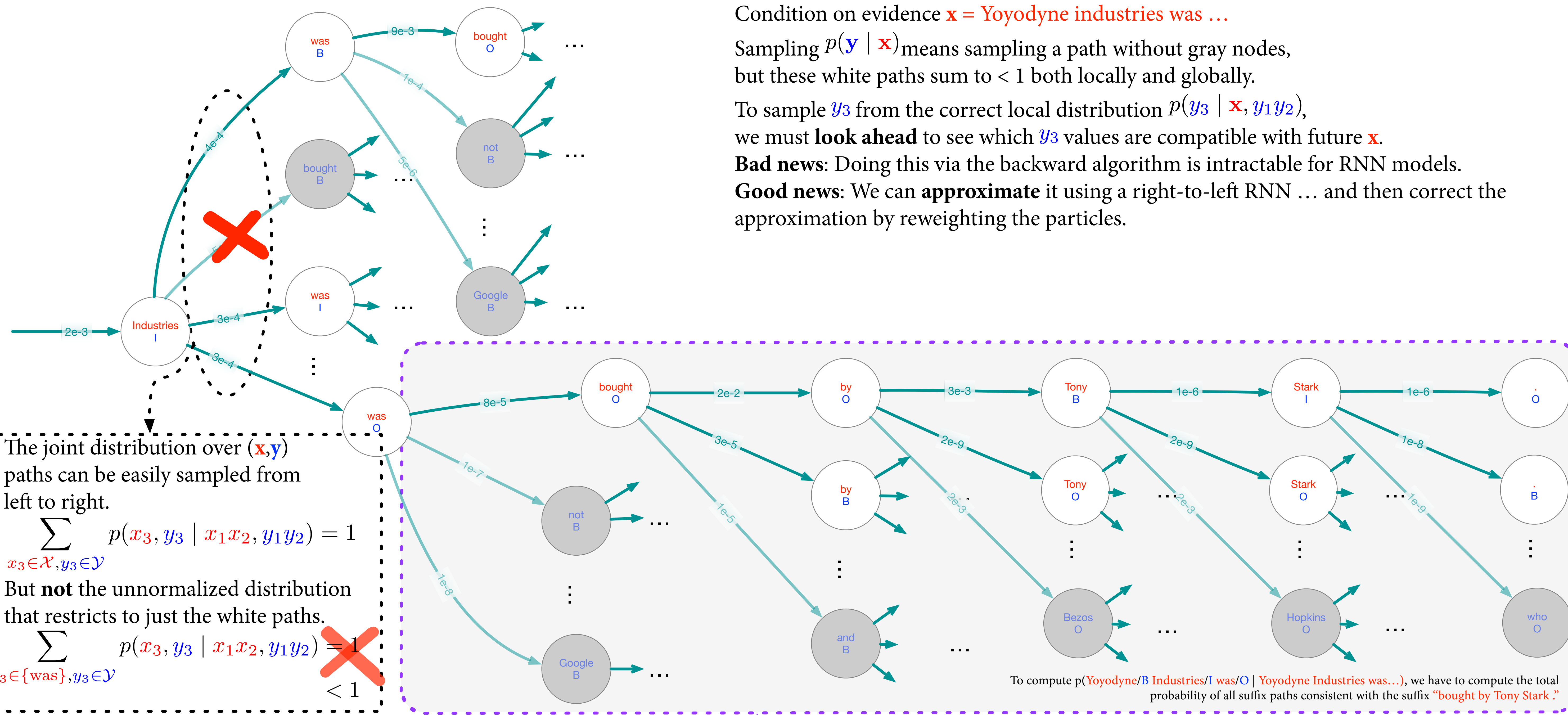the **exponentially many** $\mathbf{y}$'s. What can I do? 🤔
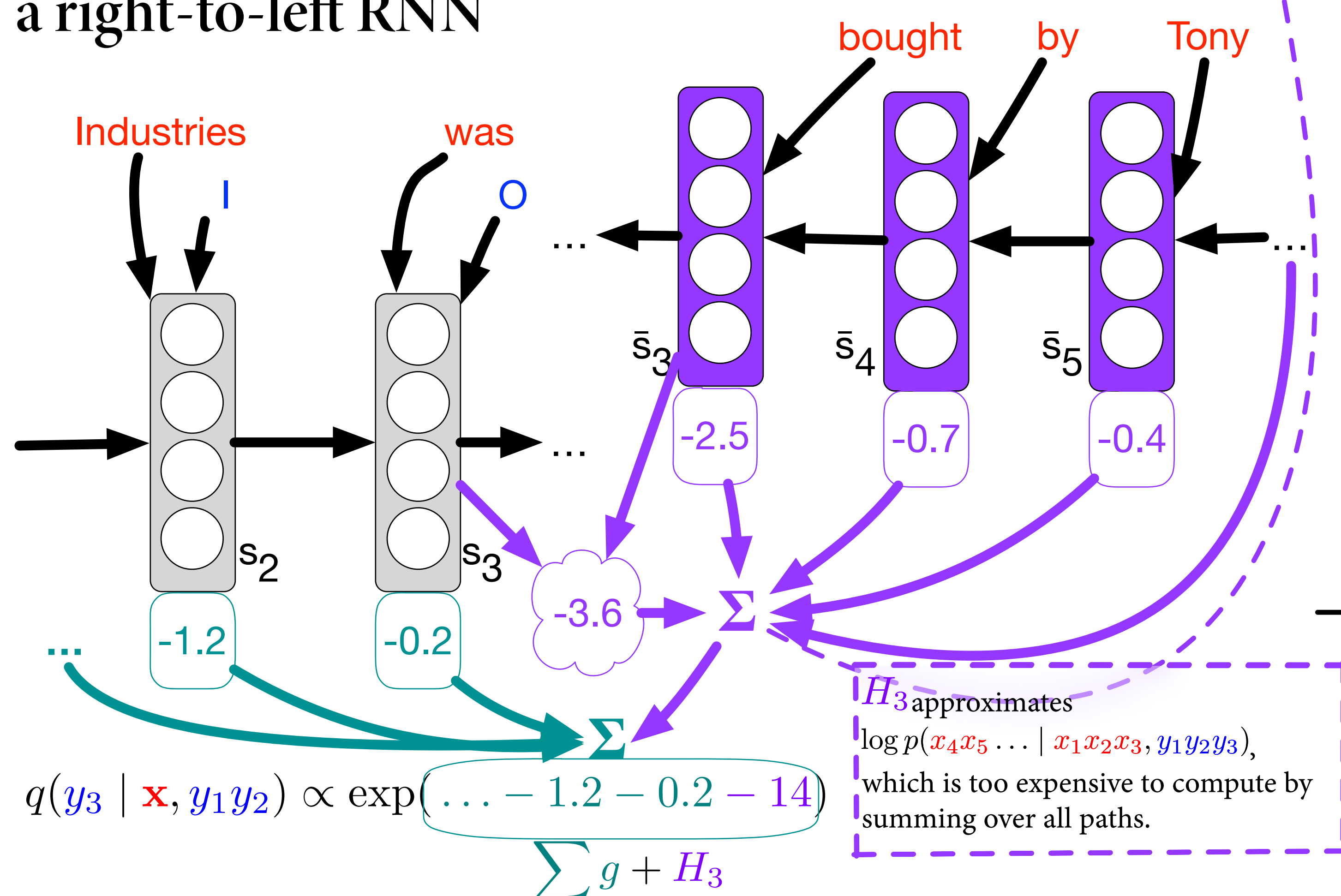
## Incremental stateful global scoring models



$\mathbf{x} =$ Yoyodyne industries was …
$\mathbf{y} =$ B I O …

$g(\mathbf{s}_{t-1}, x_t, y_t) =$ [-2.1] [-1.2] [-0.2] …

$$\log p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\left(-18.2\right)$$

e.g. $\log p(x_3, y_3 \mid x_1, y_1, x_2, y_2)$

*Can this describe a model more powerful than CRFs?*

*Could I sample "greedily" left to right? Why or why not?*

## A comparison between inference methods for sequence models

| | Left-to-right only | Left-to-right + lookahead | Data structure | Iteration step |
|---|---|---|---|---|
| Best path | Dijkstra's / Viterbi | A* | Best prefix path so far to each state | Extend a prefix chosen from a priority queue |
| Approx. best path | Beam search | Beam search + heuristic | M prefix paths of length t | Extend prefixes *exhaustively* to length t+1, then prune |
| Sampled path | Particle filtering | **Particle smoothing** | M prefix paths of length t | Extend each prefix *randomly* to length t+1, then reweight |

## Joint distribution = locally normalized distribution = greedy sampling = easy
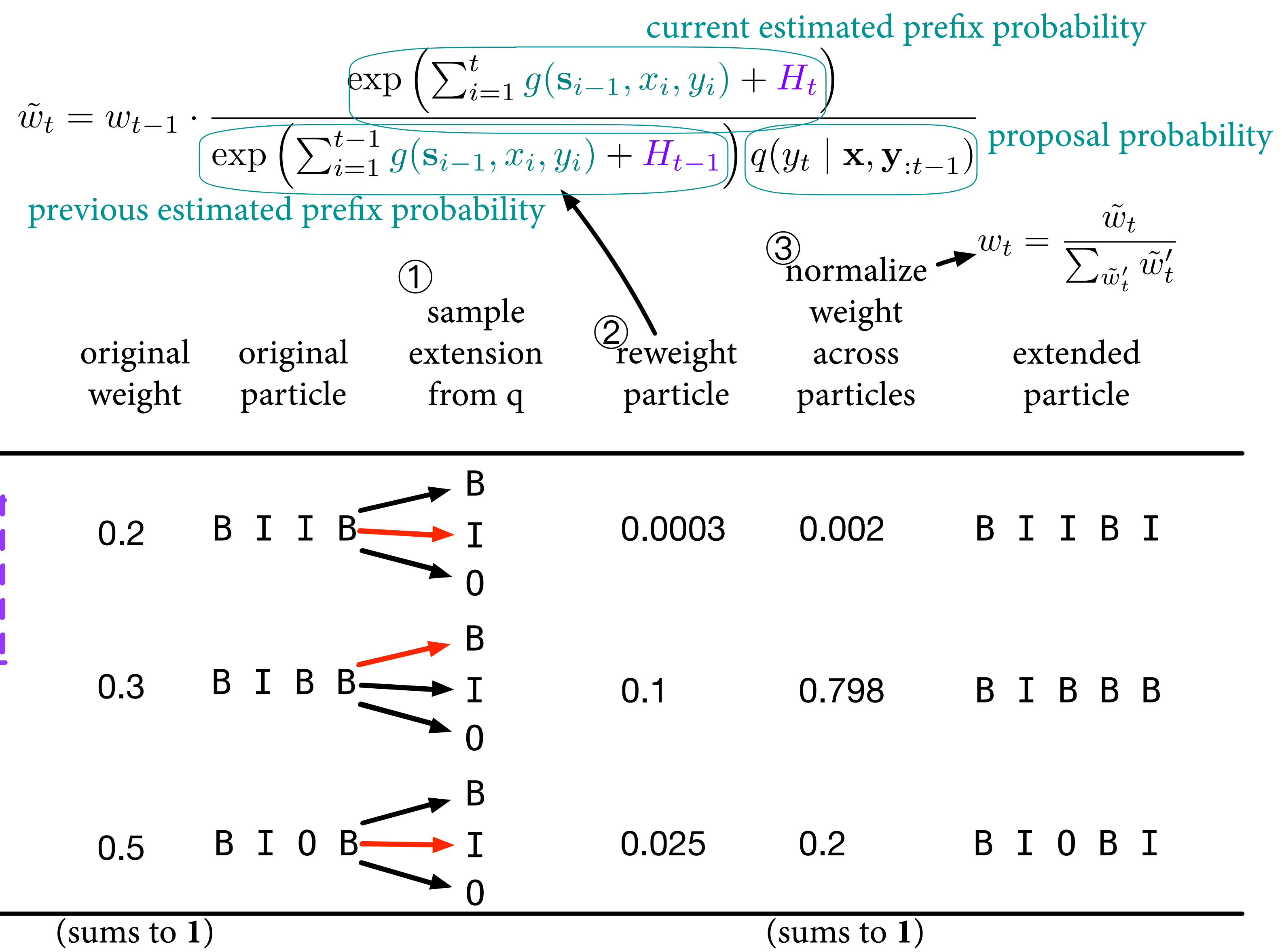## Conditional distribution = globally normalized distribution = hard!



Condition on evidence $\mathbf{x} =$ Yoyodyne industries was …

Sampling $p(\mathbf{y} \mid \mathbf{x})$ means sampling a path without gray nodes,
but these white paths sum to < 1 both locally and globally.

To sample $y_3$ from the correct local distribution $p(y_3 \mid \mathbf{x}, y_1 y_2)$,
we must **look ahead** to see which $y_3$ values are compatible with future $\mathbf{x}$.
**Bad news**: Doing this via the backward algorithm is intractable for RNN models.
**Good news**: We can **approximate** it using a right-to-left RNN … and then correct the
approximation by reweighting the particles.

The joint distribution over $(\mathbf{x}, \mathbf{y})$ paths can be easily sampled from left to right.

$$\sum_{x_3 \in \mathcal{X}, y_3 \in \mathcal{Y}} p(x_3, y_3 \mid x_1 x_2, y_1 y_2) = 1$$

But **not** the unnormalized distribution that restricts to just the white paths.

$$\sum_{x_3 \in \{was\}, y_3 \in \mathcal{Y}} p(x_3, y_3 \mid x_1 x_2, y_1 y_2) \neq 1 \quad < 1$$

To compute p(Yoyodyne/B Industries/I was/O | Yoyodyne Industries was…), we have to compute the total probability of all suffix paths consistent with the suffix "bought by Tony Stark."

## Approximating the backward algorithm with a right-to-left RNN



Industries / was    bought  by  Tony
I / O                $\bar{\mathbf{s}}_3$ $\bar{\mathbf{s}}_4$ $\bar{\mathbf{s}}_5$
[-2.5] [-0.7] [-0.4]
$\mathbf{s}_2$ [-1.2] $\mathbf{s}_3$ [-0.2] [-3.6] $\Sigma$

$q(y_3 \mid \mathbf{x}, y_1 y_2) \propto \exp\left(\dots -1.2 - 0.2 - 14\right)$
$\sum g + H_3$

$H_3$ approximates $\log p(x_4 x_5 \dots \mid x_1 x_2 x_3, y_1 y_2 y_3)$, which is too expensive to compute by summing over all paths.

We approximate $p(\mathbf{y} \mid \mathbf{x})$ with a proposal distribution
$q(\mathbf{y} \mid \mathbf{x}) = q(y_1 \mid \mathbf{x}) q(y_2 \mid \mathbf{x}, y_1) q(y_3 \mid \mathbf{x}, y_1 y_2) \dots$
We then use importance reweighting to (mostly) correct for the approximation.
This becomes accurate as the number of particles goes to infinity. We train $H_t$ to minimize the KL-divergence between q and the true distribution p:
$$\frac{\mathrm{KL}(p\|q) + \mathrm{KL}(q\|p)}{2}$$

## Sampling and weight updating steps in particle smoothing

current estimated prefix probability

$$\tilde{w}_t = w_{t-1} \cdot \frac{\exp\left(\sum_{i=1}^{t} g(\mathbf{s}_{i-1}, x_i, y_i) + H_t\right)}{\exp\left(\sum_{i=1}^{t-1} g(\mathbf{s}_{i-1}, x_i, y_i) + H_{t-1}\right) q(y_t \mid \mathbf{x}, \mathbf{y}_{:t-1})}$$

previous estimated prefix probability / proposal probability

③ normalize weight across particles $w_t = \frac{\tilde{w}_t}{\sum \tilde{w}'_t}$

| original weight | original particle | ① sample extension from q | ② reweight particle | ③ normalize weight across particles | extended particle |
|---|---|---|---|---|---|
| 0.2 | B I I B → B/I/O | 0.0003 | 0.002 | | B I I B I |
| 0.3 | B I B B → B/I/O | 0.1 | 0.798 | | B I B B B |
| 0.5 | B I O B → B/I/O | 0.025 | 0.2 | | B I O B I |
| (sums to **1**) | | | (sums to **1**) | | |

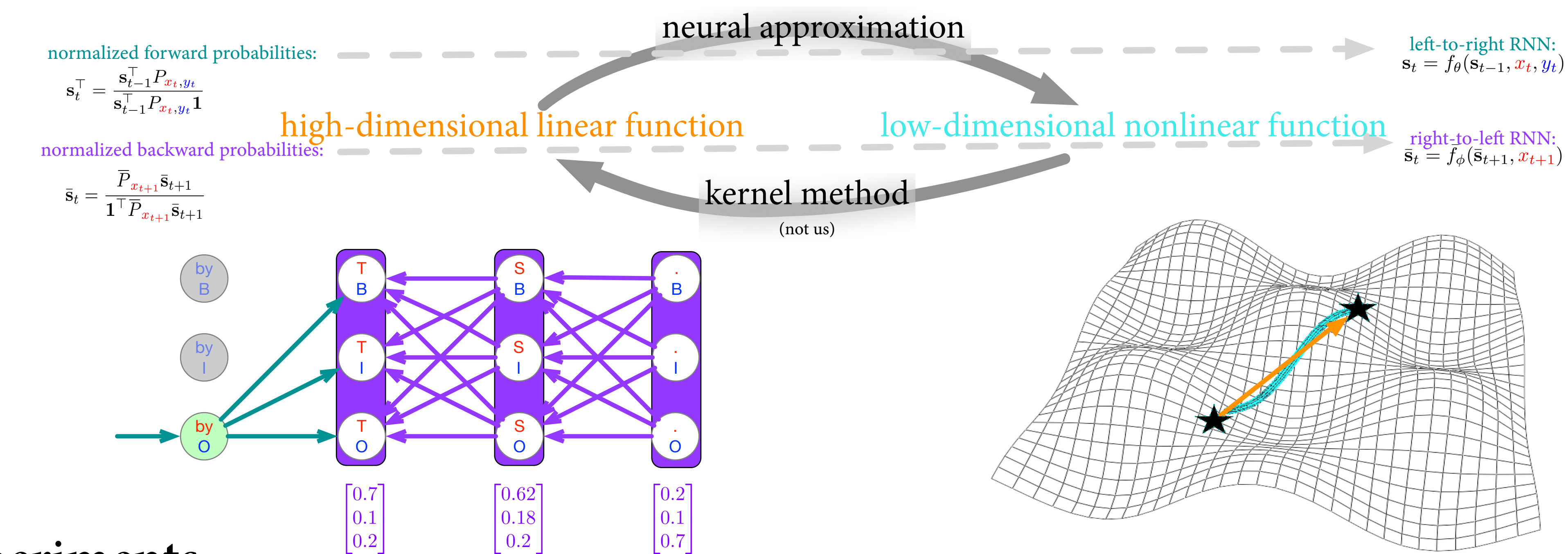## Justifying the neural approximation of dynamic programming

If the model were an HMM emitting $(x_t, y_t)$ pairs, instead of an RNN, our architecture would be exact.
The hidden vectors $\mathbf{s}_t$ and $\bar{\mathbf{s}}_t$ would represent forward and backward distributions over the hidden HMM state, and would be updated linearly at each t.
(Updates are deterministic because these are belief states, i.e., distributions over states.)
Our RNN is trying to approximate some 10000000-dimensional (?) HMM using only a 50-dimensional vector.
We hope the belief states tend to fall near a 50-dimensional manifold so that $\mathbf{s}_t$ can give the manifold coordinates. The deterministic update in that coordinate space becomes nonlinear.



normalized forward probabilities: $\mathbf{s}_t' = \frac{\mathbf{s}_{t-1}' P_{x_t, y_t}}{\mathbf{s}_{t-1}' P_{x_t, y_t} \mathbf{1}}$

normalized backward probabilities: $\bar{\mathbf{s}}_t = \frac{\bar{P}_{x_{t+1}} \bar{\mathbf{s}}_{t+1}}{\mathbf{1}^\top \bar{P}_{x_{t+1}} \bar{\mathbf{s}}_{t+1}}$

neural approximation

high-dimensional linear function / low-dimensional nonlinear function

kernel method (not us)

left-to-right RNN: $\mathbf{s}_t = f_\theta(\mathbf{s}_{t-1}, x_t, y_t)$
right-to-left RNN: $\bar{\mathbf{s}}_t = f_\phi(\bar{\mathbf{s}}_{t+1}, x_{t+1})$

## Experiments

We evaluate how good our sampler is by evaluating how similar its distribution is to the true distribution, in terms of KL-divergence. We experiment on two different model formulations, and three tasks:

— **English stressed syllable prediction**  t u ʃ ɛ t  /  0 2 0 1 0

— **Chinese social media NER**  … qu le lun dun … / … O O B I …

— **Source separation**  s ɹ ə g ɹ ɹ ə ɪ n ə ʊ ʊ l ʈ ə  /  pour fill the the glass rum  1 1 2 2 2 1 2 1 2 2 2 1 2 / 2 1 1

The particle smoothing samplers (**PS/PS:R**) consistently perform better than particle filtering baselines (**PF/PF:R**).



'best effort' lower bound of KL(q||p) — Lower value ⇒ our sampling distribution better approximates the true posterior p(y|x).

M: the # of particles sampled to approximate each distribution q(y | x).

English stressed syllables (tagging) — Chinese social media NER (tagging) — PTB (source separation) — CMUdict (source separation)