

# **Parallel sampling via autospeculation**

Thuy-Duong “June” Vuong

UC San Diego

*Simons Modern Paradigms in Generalization  
Reunion Workshop, January 20-23, 2026*

Based on joint work with

Nima Anari, Carlo Baronio, CJ Chen, Alireza Haqi, Frederic Koehler and Anqi Li

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via BERT-style marginal oracle

$$\mu(X_i = x_i | X_S = x_S)$$

\_\_\_\_\_        \_\_\_\_\_ kitchen today.

renovate      10%

clean         45%

use            45%

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle
$$\mu(X_i = x_i | X_S = x_S)$$
- Sample by any-order autoregression

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle  
 $\mu(X_i = x_i | X_S = x_S)$

## Denoising diffusion (DALL-E,...)

- $\mu : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional mean oracle  
 $f(t, x) = \mathbb{E}_{Y \sim \mu, g \sim \mathcal{N}(0, \frac{1}{t} \cdot I)} [Y | Y + g = \frac{x}{t}]$

\_\_\_\_\_.

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle  
 $\mu(X_i = x_i | X_S = x_S)$

_____	
_____	
_____	
_____	
_____	
	today 15%
	Tuesday 5%
	never 20%
	....

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle  
 $\mu(X_i = x_i | X_S = x_S)$

\_\_\_\_\_ today.

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle
$$\mu(X_i = x_i | X_S = x_S)$$

Tom \_\_\_\_\_ today.

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle  
 $\mu(X_i = x_i | X_S = x_S)$

Tom will \_\_\_\_\_ today.



# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle  
 $\mu(X_i = x_i | X_S = x_S)$

Tom will eat \_\_\_\_\_ today.

# Generative modeling

## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle
$$\mu(X_i = x_i | X_S = x_S)$$

Tom will eat breakfast today.

# Generative modeling

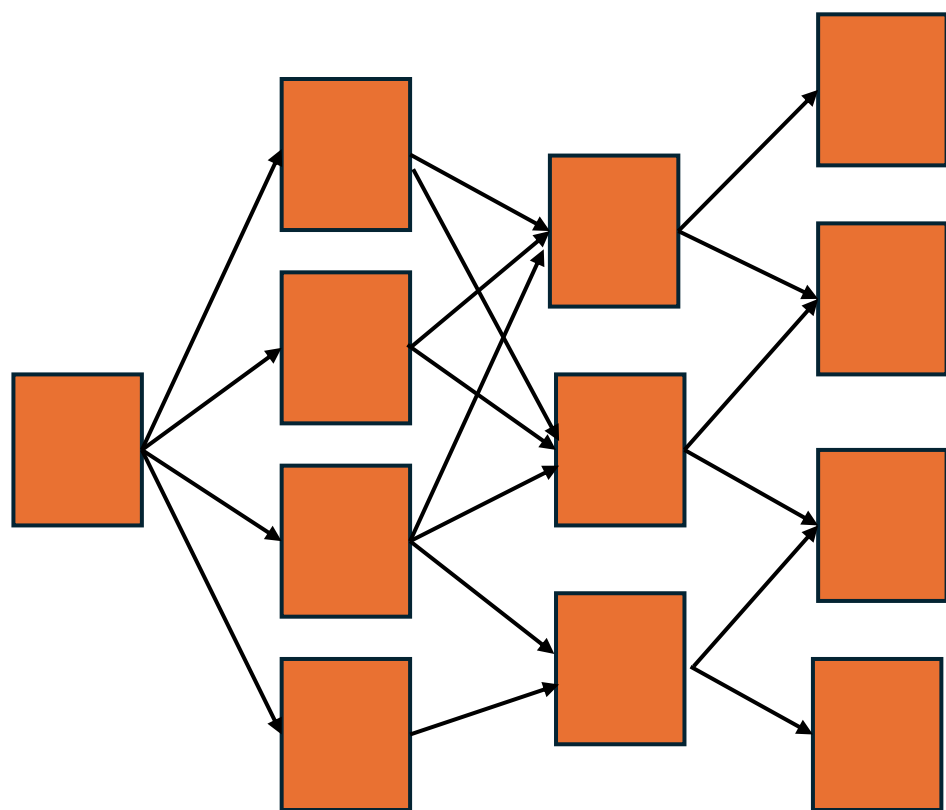
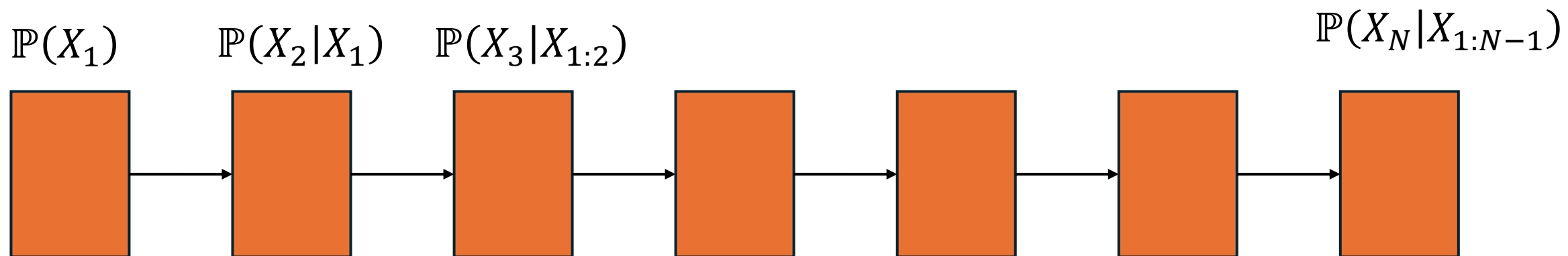
## Autoregression (ChatGPT,...)

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional marginal oracle
$$\mu(X_i = x_i | X_S = x_S)$$
- Sample by any-order autoregression

## Denoising diffusion (DALL-E,...)

- $\mu : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$
- Sample via conditional mean oracle
$$f(t, x) = \mathbb{E}_{Y \sim \mu, g \sim \mathcal{N}(0, \frac{1}{t} \cdot I)} [Y | Y + g = \frac{x}{t}]$$
- Sample by stochastic differential eq (SDE)

# Parallelization?



- Seems inherently sequential. Parallelization?
- Parallel computing model:
  - Many parallel processes per round
  - Communication at the end of each round
- Our results: autoregression (& diffusion) in  $\tilde{O}(n^{\frac{1}{2}})$  rounds and  $O(n \log n)$  total work

# Autoregression (ChatGPT,...)

General  $\mu: [q]^n \rightarrow \mathbb{R}_{\geq 0}$ , BERT marginals  $\mu(X_i = x_i | X_S = x_S)$

- Folklore:  $\Theta(n)$  oracle queries/total work
- [AGR24] :
  - Upper bound:  $\tilde{O}(n^{\frac{2}{3}})$  rounds &  $O(n \log n)$  queries/total work
  - Lower bound:  $\Omega(n^{\frac{1}{3}})$  rounds  $\forall$  poly-queries algorithm
- This work:
  - $\tilde{O}(n^{\frac{1}{2}})$  rounds &  $O(n \log n)$  queries/total work
  - Simpler algorithm that also works for diffusion

# Autoregression (ChatGPT,...)

General  $\mu: [q]^n \rightarrow \mathbb{R}_{\geq 0}$ , **BERT** marginals  $\mu(X_i = x_i | X_S = x_S)$

- Folklore:  $\Theta(n)$  oracle queries/total work
- [AGR24] :
  - Upper bound:  $\tilde{O}(n^{\frac{2}{3}})$  rounds &  $O(n \log n)$  queries/total work
  - Lower bound:  $\Omega(n^{\frac{1}{3}})$  rounds  $\forall$  poly-queries algorithm
- This work:
  - $\tilde{O}(n^{\frac{1}{2}})$  rounds &  $O(n \log n)$  queries/total work
- **BERT** marginals are **necessary** for **parallelization**:
  - **Can't** get  $o(n)$  parallel time with only **fixed-order marginals**  $\mu(X_i | X_{1:i-1})$

Algorithm

# Rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\frac{\mu}{\nu}(\mathbf{x})$

$$\rightarrow \mathbb{P}[\text{output } \mathbf{x}] = \nu(\mathbf{x}) \cdot \frac{\mu}{\nu}(\mathbf{x}) = \mu(\mathbf{x})$$



# Rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\frac{\mu}{\nu}(\mathbf{x})$  What if  $\frac{\mu}{\nu}(\mathbf{x}) > 1$  ?

$$\rightarrow \mathbb{P}[\text{output } \mathbf{x}] = \nu(\mathbf{x}) \cdot \frac{\mu}{\nu}(\mathbf{x}) = \mu(\mathbf{x})$$

# Speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$

$$\rightarrow \mathbb{P}[\text{output } \mathbf{x}] = \min\{\mu(\mathbf{x}), \nu(\mathbf{x})\}$$

**Wrong output distribution!**

# Speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$   
 $\rightarrow \mathbb{P}[\text{output } \mathbf{x}] = \min\{\mu(\mathbf{x}), \nu(\mathbf{x})\}$
3. While not accepted do:  
    Sample  $\mathbf{y} \sim \mu$   
    Accept and return  $\mathbf{y}$  with probability  $\max\{0, 1 - \frac{\nu}{\mu}(\mathbf{y})\}$   
 $\rightarrow \mathbb{P}[\text{output } \mathbf{y}] = \max\{0, \mu(\mathbf{y}) - \nu(\mathbf{y})\}$

# Recursive speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ : **How to implement?**

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$
3. While not accepted do:  
**Sample  $\mathbf{y} \sim \mu$  by:  $y_{1:n/2} \sim \mu(X_{1:n/2}), y_{n/2+1:n} \sim \mu(X_{n/2+1:n} | X_{1:n/2} = y_{1:n/2})$**   
Accept and return  $\mathbf{y}$  with probability  $\max\{0, 1 - \frac{\nu}{\mu}(\mathbf{y})\}$

# Recursive speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution that admits fast sampler e.g. in  $O(1)$  rounds  
& can be built w/ conditional marginal/mean oracles for  $\mu$

The following exactly samples from  $\mu$ : **How to implement?**

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$
3. **While not accepted do:**  
...

- Fully sequential: no gain
- Fully parallel: need infinite #processors

# Recursive speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :      **How to implement?**

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$

3. **While not accepted do:**

...

- Fully sequential: no gain
- Fully parallel: need infinite #processors
- Key idea: **sequentially** process batches of geometrically increasing size, where each batch is processed **in parallel**

# Recursive speculative rejection sampling

$\mu$ : target distribution

$\nu$ : reference distribution

The following exactly samples from  $\mu$ :

1. Sample  $\mathbf{x} \sim \nu$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu}{\nu}(\mathbf{x})\}$
3. For  $r=0,1,2,\dots$  do:  
    For  $i = \lceil (1 + \rho)^r \rceil, \lceil (1 + \rho)^r \rceil + 1, \dots, \lceil (1 + \rho)^{r+1} \rceil - 1$  do in parallel  
    Sample  $\mathbf{y} \sim \mu$  by:  
         $y_{1:n/2} \sim \mu(X_{1:n/2}), y_{n/2+1:n} \sim \mu(X_{n/2+1:n} | X_{1:n/2} = y_{1:n/2})$   
    Accept and return  $\mathbf{y}$  with probability  $\max\{0, 1 - \frac{\nu}{\mu}(\mathbf{y})\}$

# Recursive speculative rejection sampling

- $S = \{a + 1, \dots, b\}$
- $\mu_S$ : target distribution
- $\nu_S$ : reference distribution

In  $O(1)$  rounds and  $O(|S|)$  queries, can:

- Compute  $\frac{d\mu_S}{d\nu_S}$
- Sample  $\nu_S$

The following exactly samples from  $\mu_S$ :

1. Sample  $\mathbf{x} \sim \nu_S$
2. Accept and return  $\mathbf{x}$  with probability  $\min\{1, \frac{\mu_S}{\nu_S}(\mathbf{x})\}$
3. For  $r=0,1,2,\dots$  do:
  - For  $i = \lceil (1 + \rho)^r \rceil, \lceil (1 + \rho)^r \rceil + 1, \dots, \lceil (1 + \rho)^{r+1} \rceil - 1$  do in parallel:
  - Sample  $\mathbf{y} \sim \mu_S$  by: Set  $L(S) = \{a + 1, \dots, \frac{a+b}{2}\}$ ,  $R(S) = \{\frac{a+b}{2} + 1, \dots, b\}$
  - $\mathcal{Y}_{L(S)} \sim \mu_{L(S)}, \mathcal{Y}_{R(S)} \sim \mu_{R(S)}(\cdot | X_{<R(S)} = \dots || \mathcal{Y}_{L(S)})$
  - Accept and return  $\mathbf{y}$  with probability  $\max\{0, 1 - \frac{\nu_S}{\mu_S}(\mathbf{y})\}$

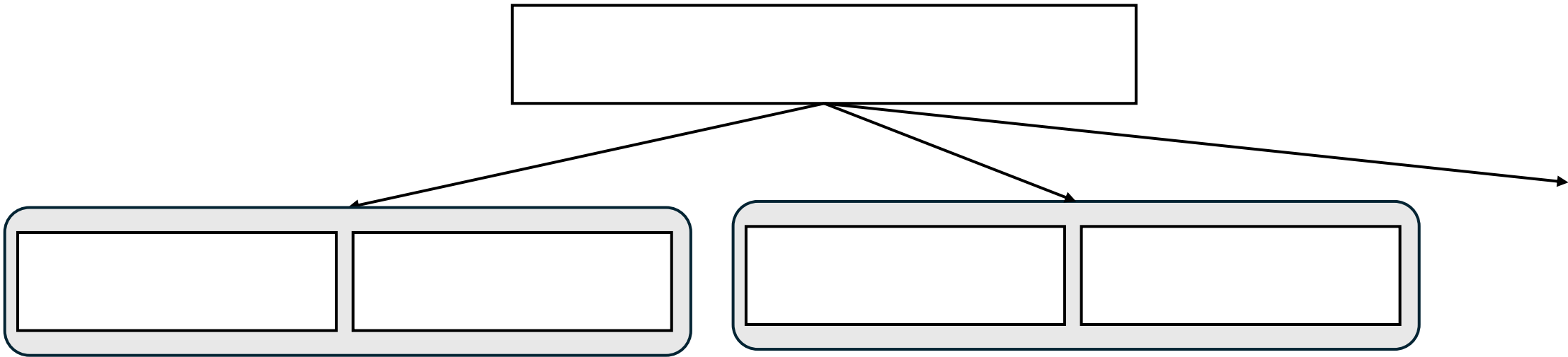


# What are $\mu_S$ and $\nu_S$ ?

- $\mu : [q]^n \rightarrow \mathbb{R}_{\geq 0}$
  - Oracle access to BERT-style marginals  $\mu(X_i = x_i | X_S = x_S)$
  - Sample order  $\sigma \sim \text{Uniform}(S_n)$ 
    - $S = \{a + 1, \dots, b\}$
    - $\mu_S = \mu(X_{\sigma(S)} | X_{\sigma(<S)}) = \mu(X_{\sigma(a+1:b)} | X_{\sigma(1:a)})$
    - $\nu_S$  = product distribution with same marginal as  $\mu_S$   
 $= \mu(X_{\sigma(a+1)} | X_{\sigma(1:a)}) \otimes \dots \otimes \mu(X_{\sigma(b)} | X_{\sigma(1:a)})$   
Can sample each coordinate of  $\nu_S$  in parallel  
b/c no dependencies
- For fixed order  $\sigma$ ,  
 $TV(\mu_S, \nu_S) = \Omega(1)$   
b/c dependencies
- Pinning lemma:  
 $\mathbb{E}_{\sigma \sim S_n}[TV(\mu_S, \nu_S)]$   
 $\approx |S|/\sqrt{n}$

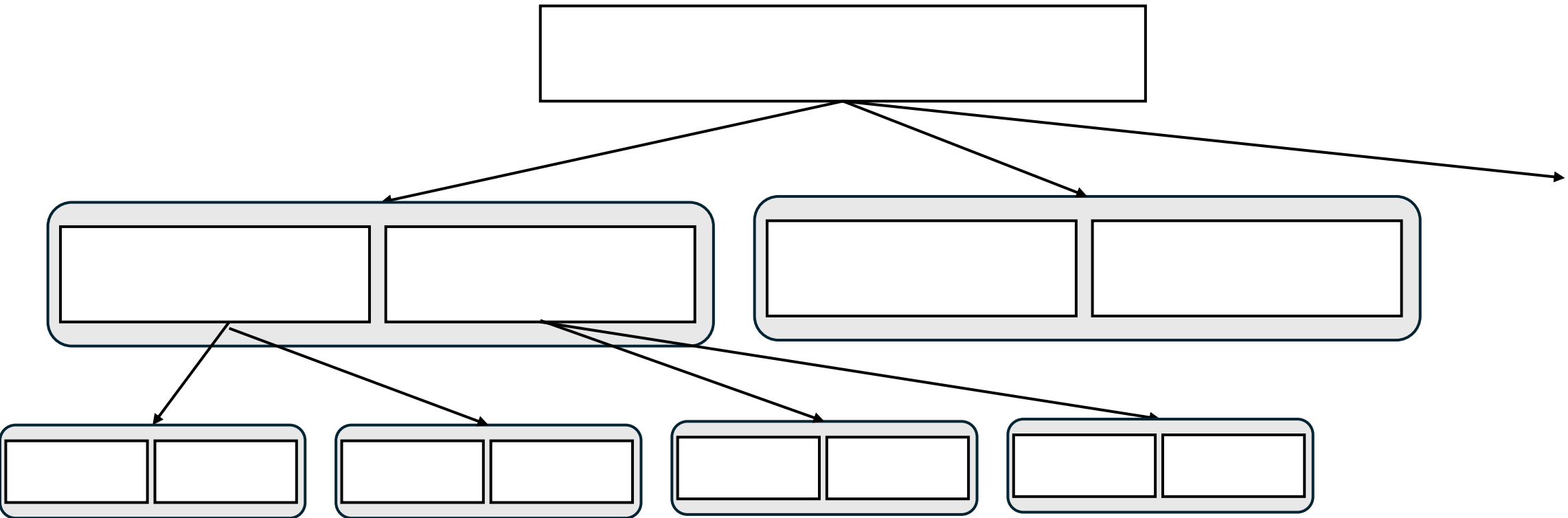
# Analysis: high level

- Difficulty: Many parallel threads & deep recursion



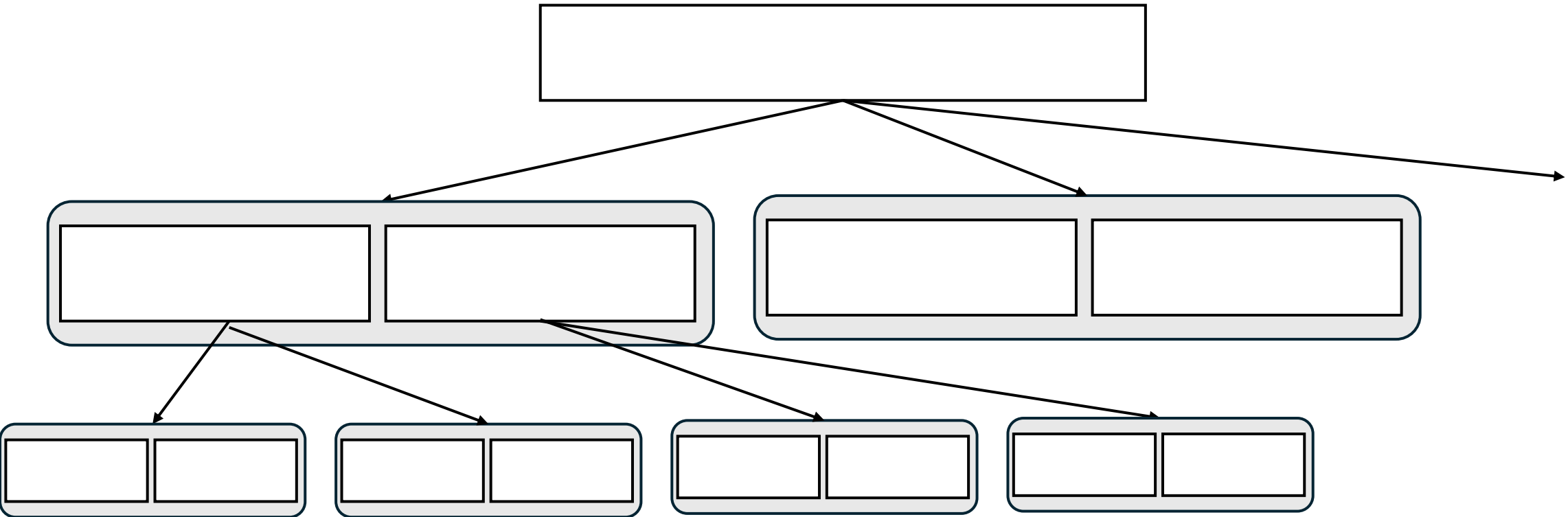
# Analysis: high level

- Difficulty: Many parallel threads & deep recursion



# Analysis: high level

- Difficulty: Many parallel threads & deep recursion
- $\mathbb{E}[T] = \mathbb{E}[\sum_S TV(\mu_S, \nu_S)] \approx \tilde{O}(\sqrt{n})$  **Pinning lemma**



# Summary & open questions

General  $\mu: [q]^n \rightarrow \mathbb{R}_{\geq 0}$ , BERT marginals  $\mu(X_i = x_i | X_S = x_S)$

- [AGR24] :
  - Upper bound:  $\tilde{O}(n^{\frac{2}{3}})$  rounds &  $O(n \log n)$  queries/total work
  - Lower bound:  $\Omega(n^{\frac{1}{3}})$  rounds  $\forall$  poly-queries algorithm
- This work:
  - $\tilde{O}(n^{\frac{1}{2}})$  rounds &  $O(n \log n)$  queries/total work
  - Simpler algorithm that also works for diffusion
- Open:
  - Handling more noise via (parallelized?) back-tracking