

# Adaptive Data Collection via Autoregressive Generation

**Hong Namkoong**

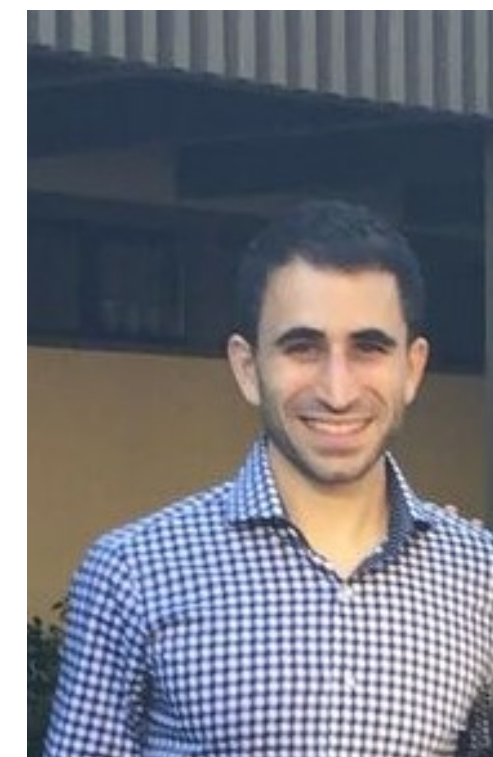
Columbia University



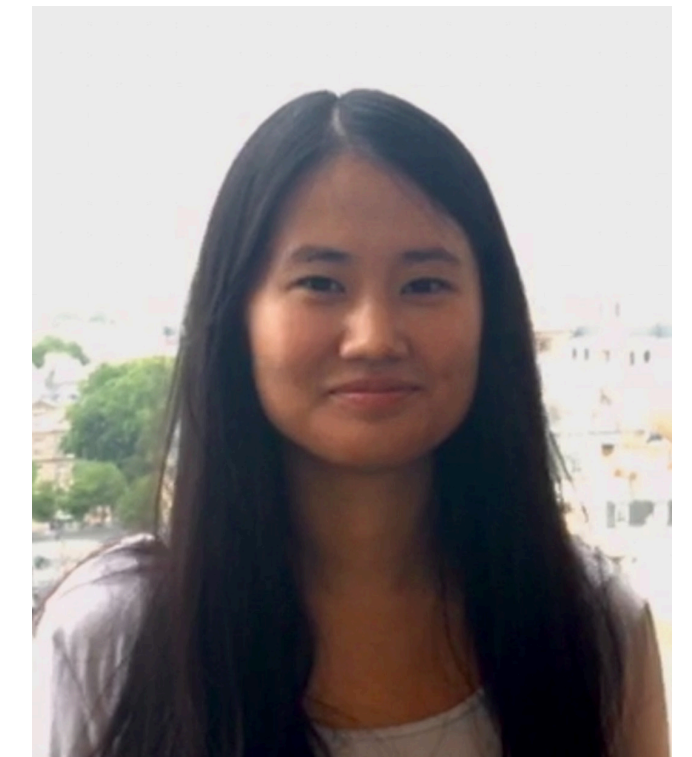
Naimeng Ye



Tiffany Cai



Dan Russo



Kelly Zhang

# My journey on distribution shift

- Algo for specific shifts => scale data => what if web-data isn't enough?

# My journey on distribution shift

- Algo for specific shifts => scale data => what if web-data isn't enough?

## Intellectual bottleneck: no language for datasets

Inductive modeling for distribution shifts, esp.  $Y|X$ -shifts

Liu, Wang, Cui, N. On the need for a language describing distribution shifts: Illustrations on tabular datasets. NeurIPS23.

Cai, N., Yadlowsky. Diagnosing model performance under distribution shift. FORC '23, journal version under revision in Operations Research.

NeurIPS23 tutorial: <https://nips.cc/virtual/2023/tutorial/73953>

# Today: Uncertainty

1. Understand how different current dataset is from training
2. Adaptively collect data on distributions that can reduce uncertainty

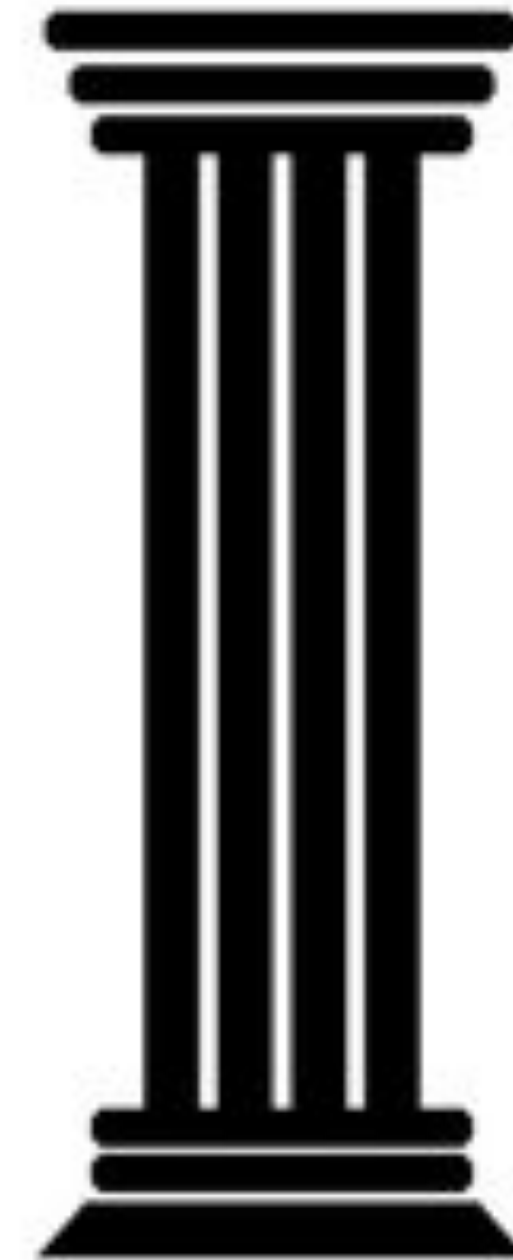


# Uncertainty

- Intelligent agents must comprehend uncertainty and take actions to resolve it
- Several line of work tackle this problem
  - Bayesian neural networks, GPs, ensembles, epistemic neural nets, conformal prediction, multi-calibration...many other interesting ideas
- But these ideas have not materialized in the form of scalable models
  - In the sense that they are not incorporated into Llama3

# Why?

- Two pillars of ML
  1. Optimize fictitious loss on web-scale data
  2. Test engineering innovations based on val loss
- Hard to fit aforementioned ideas into this umbrella



**Today:** Adopt these principles to quantify uncertainty

# Classical Approach

Model “environment” first, then pass onto quantity of interest

- Bayes rule provides a natural modeling language

Latent “environment” drawn  $U \sim \pi(\cdot)$

Data generated by environment  $\mathbb{P}(Y | X, U)$

- As you gather data from an environment, infer what the environment looks like

Posterior  $\mathbb{P}(U | \text{Data})$

# Example: mental disease diagnosis

## Why probabilistic modeling is hard

- Goal: uncertainty quantification on diagnoses
  - Prob ( schizophrenia | Q & A with patient )
- Latent parameter: patient's "mental health state"



# Example: mental disease diagnosis

## Why probabilistic modeling is hard

- Goal: uncertainty quantification on diagnoses
  - Prob ( schizophrenia | Q & A with patient )
- Latent parameter: patient's "mental health state"

**Latent has no physical meaning!**



# Example: mental disease diagnosis

## Why probabilistic modeling is hard

- Goal: uncertainty quantification on diagnoses
  - Prob ( schizophrenia | Q & A with patient )
- Latent parameter: patient's "mental health state"

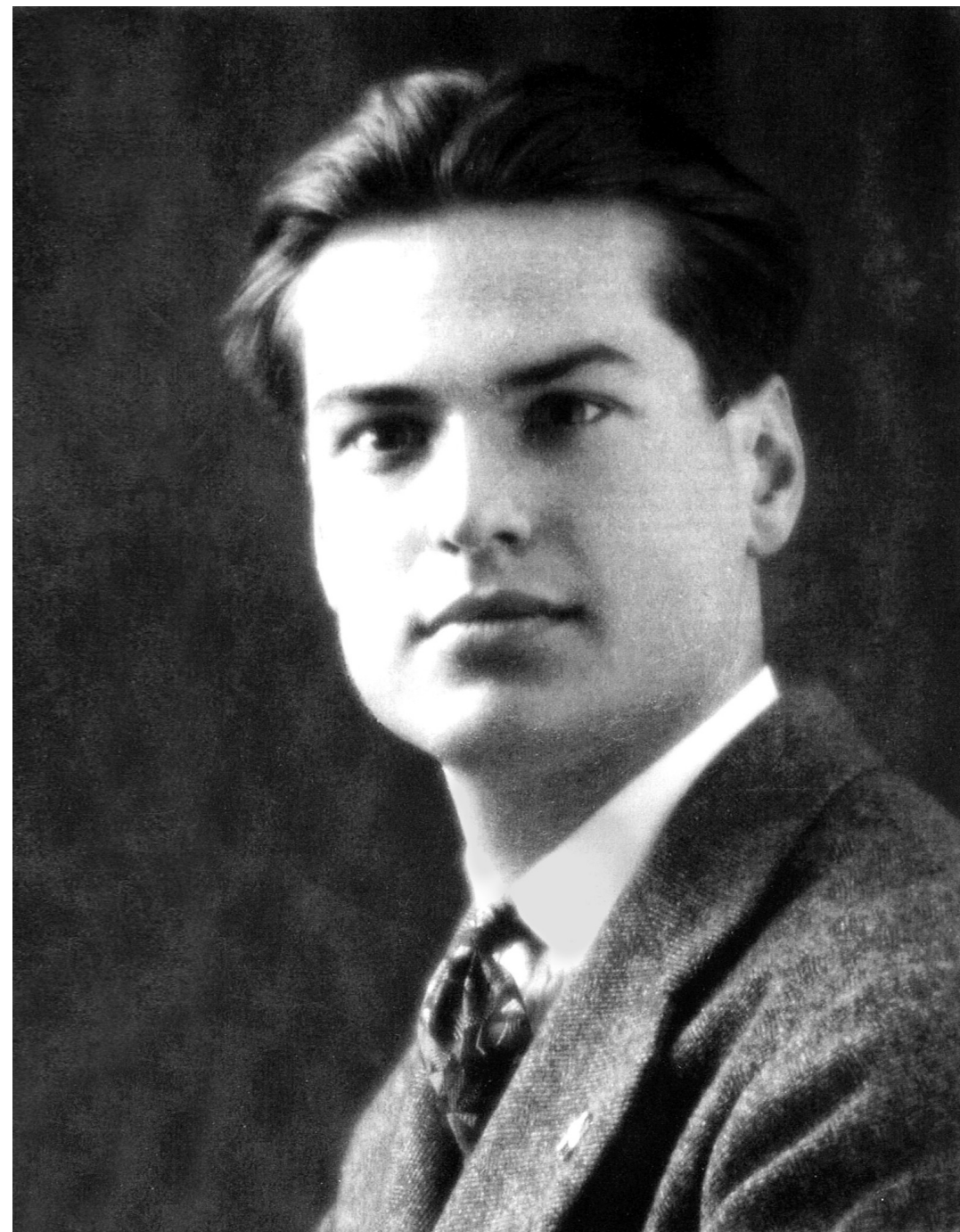
**Latent has no physical meaning!**

**Hard to check whether your unicorn is better than mine**



# Uncertainty comes from missing data

## Bayesian modeling a la De Finetti



Bruno De Finetti  
(1929)

“If data exchangeable, autoregressive modeling  
= modeling latent environment”

Modeling primitive: autoregressive probability

$\mathbb{P}$  (next data | past data, patient)

# Probabilistic modeling as pre-training

## Modern interpretation of De Finetti's philosophy

- Utilize **massive** historical data across many environments
- Modeling primitive: given past observations, prob of next observation

$$\text{sequence prediction loss} = \sum_{(\text{patients, data})} \log \hat{P} (\text{next data} \mid \text{past data, patient})$$

# Probabilistic modeling as pre-training

## Modern interpretation of De Finetti's philosophy

- Utilize **massive** historical data across many environments
- Modeling primitive: given past observations, prob of next observation

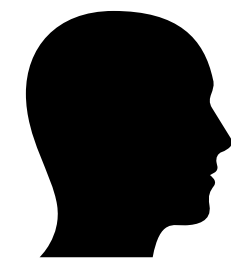
$$\text{sequence prediction loss} = \sum_{(\text{patients, data})} \log \hat{P}(\text{next data} \mid \text{past data, patient})$$

- Great news: we know how to do sequence modeling well!
  - Above measure is exactly what we build scaling laws on



# Conceptual example with language: mental health

## Diagnosing based on verbal sessions

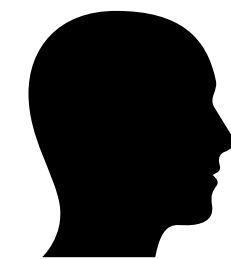


How have you been feeling emotionally over the past few weeks?

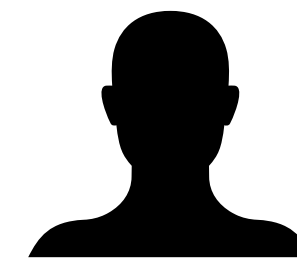


I've been feeling really overwhelmed lately. It's like I can't keep up with everything...

⋮

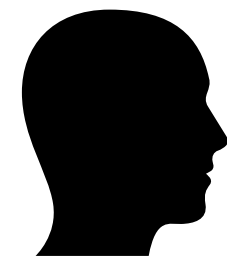


How about your sleep and appetite? Have you noticed any changes there?

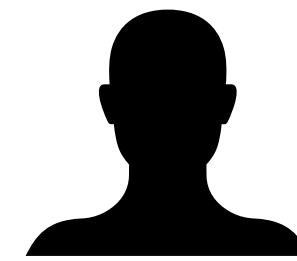


I've been sleeping a lot, but it doesn't feel restful. And I haven't had much of an...

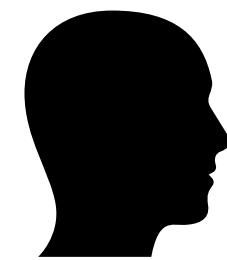
# Simulation 1



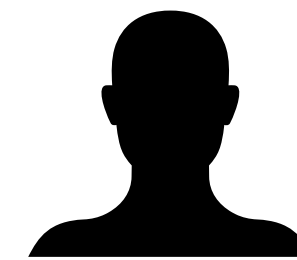
How have you been feeling emotionally over the past few weeks?



I've been feeling really overwhelmed lately. It's like I can't keep up with everything...



How about your sleep and appetite? Have you noticed any changes there?



I've been sleeping a lot, but it doesn't feel restful. And I haven't had much of an...



Have you noticed any changes in your thoughts or perceptions recently? For example, have you ever seen or heard things that others didn't seem to notice?



Sometimes I hear voices. They're not always clear, but I can hear them talking, even when no one's around.

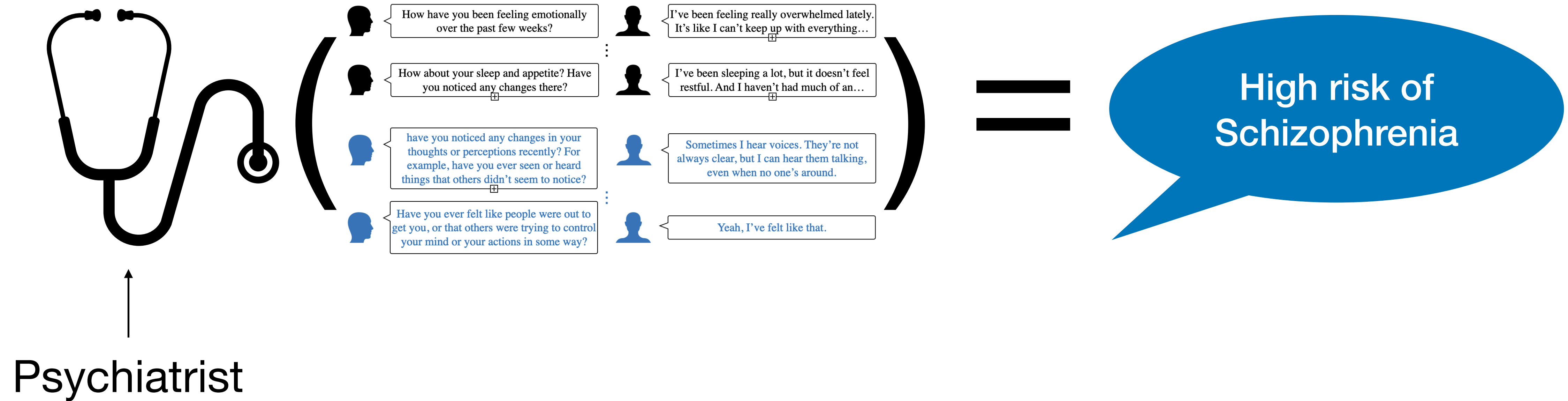


Have you ever felt like people were out to get you, or that others were trying to control your mind or your actions in some way?



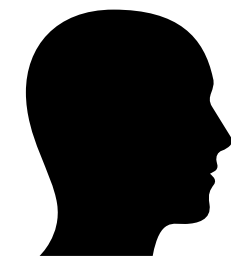
Yeah, I've felt like that.

# Simulation 1

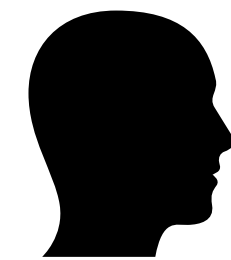




# Simulation 1



How have you been feeling emotionally over the past few weeks?



How about your sleep and appetite? Have you noticed any changes there?



have you noticed any changes in your thoughts or perceptions recently? For example, have you ever seen or heard things that others didn't seem to notice?



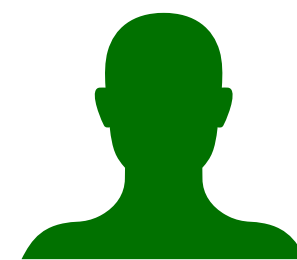
Have you ever felt like people were out to get you, or that others were trying to control your mind or your actions in some way?



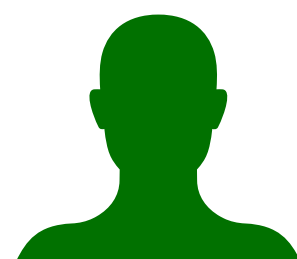
I've been feeling really overwhelmed lately. It's like I can't keep up with everything...



I've been sleeping a lot, but it doesn't feel restful. And I haven't had much of an...

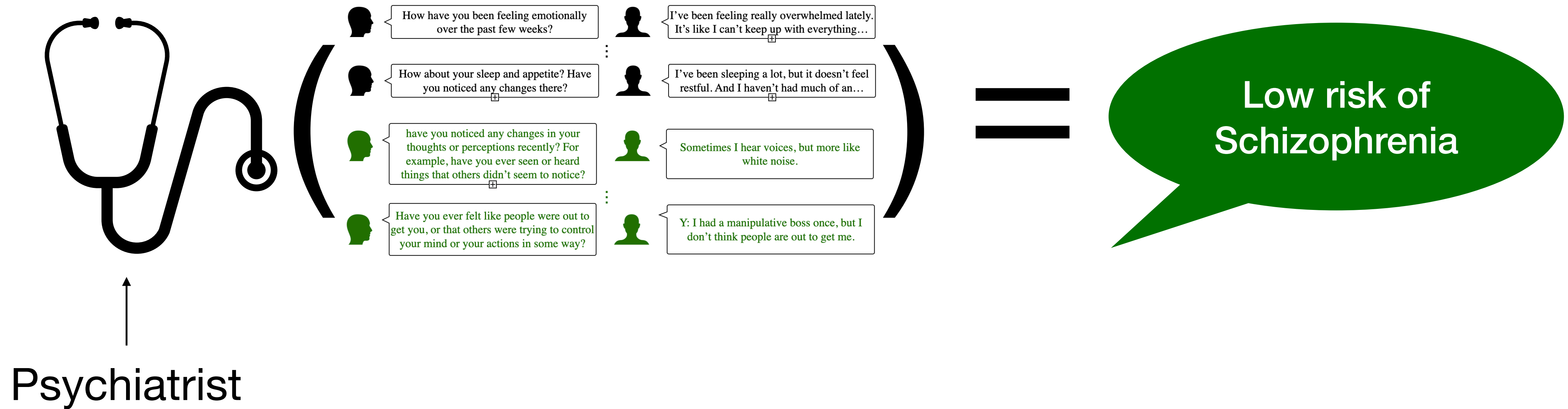


Sometimes I hear voices, but more like white noise.



Y: I had a manipulative boss once, but I don't think people are out to get me.

# Simulation 1



# Conceptual example: diagnosis based on verbal sessions

Observed

X: How have you been feeling emotionally over the past few weeks?

⋮

X: How about your sleep and appetite? Have you noticed any changes there?

Y: Feeling really overwhelmed. It's been affecting my mood a lot.

⋮

Y: Sleeping a lot, but doesn't feel restful. No appetite, some days I forget to eat.

Generated

X: Have you noticed any changes in your thoughts or perceptions recently? For example, have you ever seen or heard things that others didn't seem to notice?

⋮

X: Have you ever felt like people were out to get you, or that others were trying to control your mind or your actions in some way?

Y: Sometimes I hear voices. They're not always clear, but I can hear them talking, even when no one's around.

⋮

Y: I've often felt like that. Specifically, I'm worried my parents are conspiring to steal from me.

Y: Sometimes I hear voices, but it's more like white noise. Is that what you meant?

⋮

Y: I had a manipulative boss once, but I don't think people are out to get me.



High risk of schizophrenia

...



Low risk of schizophrenia

**Main insight:**  
variability in inferred state across s = uncertainty in diagnosis

# Conceptual example: diagnosis based on verbal sessions

Observed

X: How have you been feeling emotionally over the past few weeks?

⋮

X: How about your sleep and appetite? Have you noticed any changes there?

Y: It's been a rollercoaster. Some days, my mind feels like it's working against me.

⋮

Y: I have vivid nightmares, and have great difficulty sleeping.

Generated

X: Have you noticed any changes in your thoughts or perceptions recently? For example, have you ever seen or heard things that others didn't seem to notice?

⋮

X: Have you ever felt like people were out to get you, or that others were trying to control your mind or your actions in some way?

Y: Sometimes I hear voices. They're not always clear, but I can hear them talking, even when no one's around.

⋮

Y: I've often felt like that. Specifically, I'm worried my parents are conspiring to steal from me.

Y: I see shadows of animals chasing me, and sometimes hear a voice saying "you're not good enough"

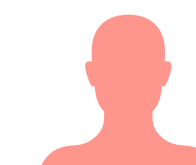
⋮

Y: I'm worried my co-worker is trying to sabotage my career.



High risk of schizophrenia

...



High risk of schizophrenia

**Main insight:**

variability in inferred state across s = uncertainty in diagnosis

# Autoregressive vs. marginal predictions

## Confusion #1 in the literature

- Autoregressive generation critical
  - Averaging future data washes out idiosyncratic/aleatoric uncertainty
  - Remaining correlation reflect epistemic uncertainty
- Current ML literature on probabilistic interpretations of sequence learning do not differentiate between epistemic vs. aleatoric uncertainty

# Related work

## **TLDR: correctness of autoregressive generation known since 1920s**

- De Finetti [1929] showed modeling of exchangeable sequences of observable RVs is equivalent to Bayesian modeling of latent parameters
- Bayesian multiple imputation views of casual inference: Rubin [1978]
- Related ideas (re)discovered and articulated many times in many communities
  - Math and philosophy of exchangeable sequence modeling: Berti and coauthors [1998, 2021, 2022], Fortini et al. [2014, 2023], Fong, Holmes, and Walker [2023]
  - Neural processes: Garnelo et. al [2203]
  - Joint predictions: Osband et. al [2022]

# So what? AI-driven decisions

- AI now comprehends language and visual inputs
- Big opportunities to make decisions based on them
- Decision-making requires comprehending uncertainty and acting to resolve it



Cold start problem in RecSys

# So what? AI-driven decisions

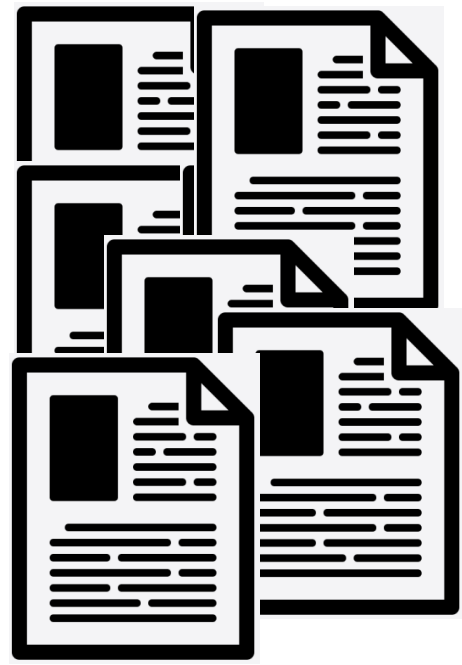
- AI now comprehends language and visual inputs
- Big opportunities to make decisions based on them
- Decision-making requires comprehending uncertainty and acting to resolve it





# Today: news recommendations

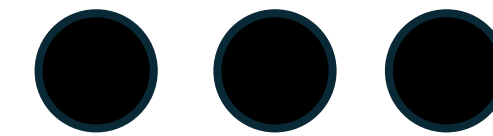
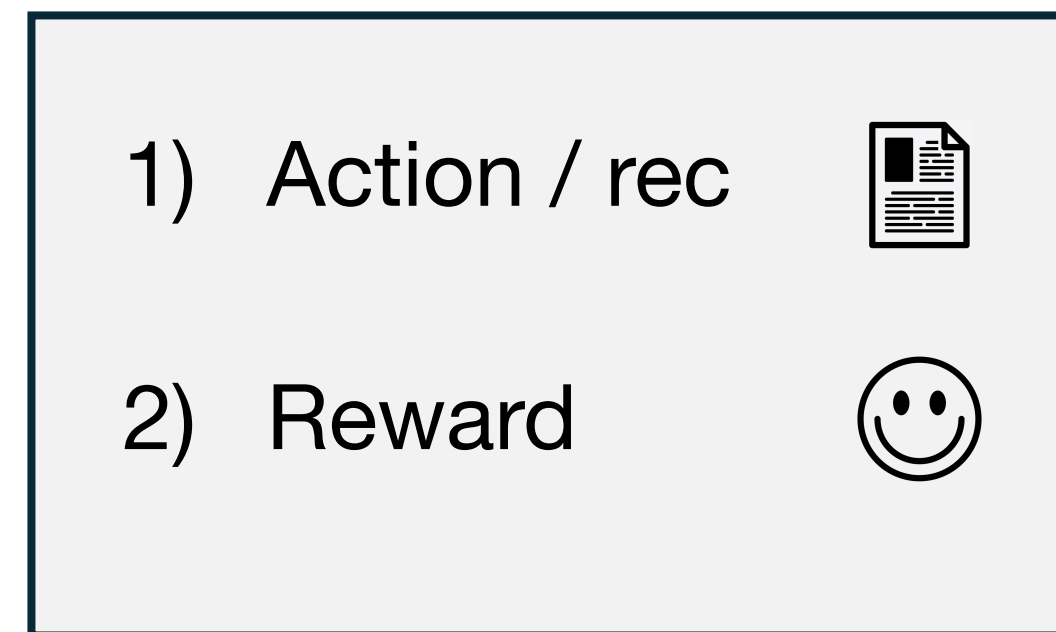
New articles  
are released



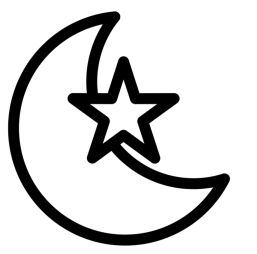
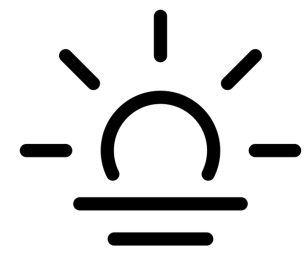
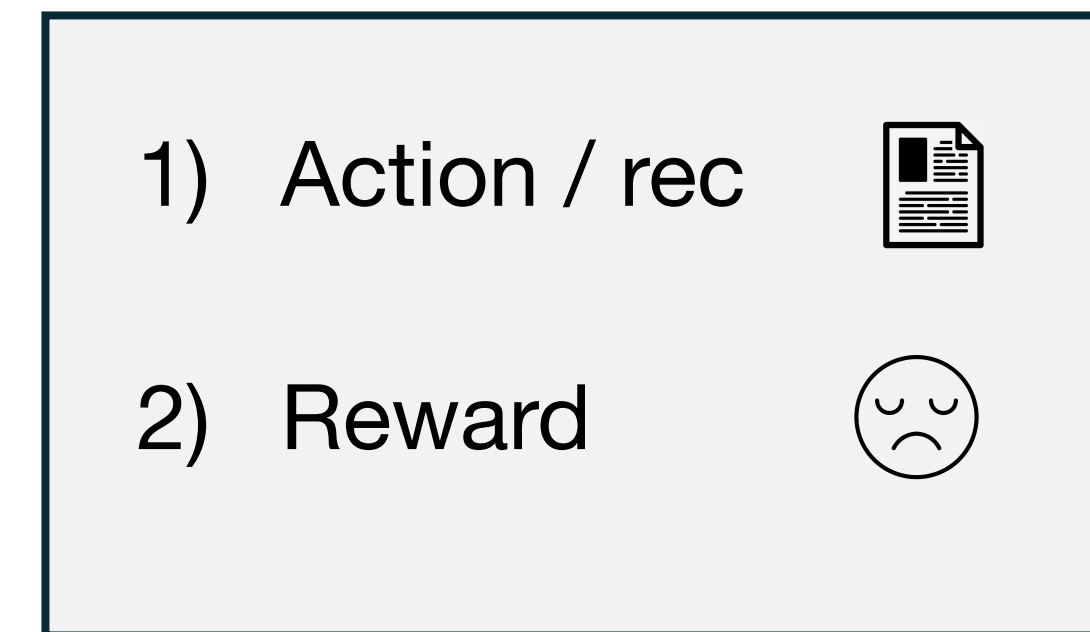
An LLM  
reads them



 Interact with User 1

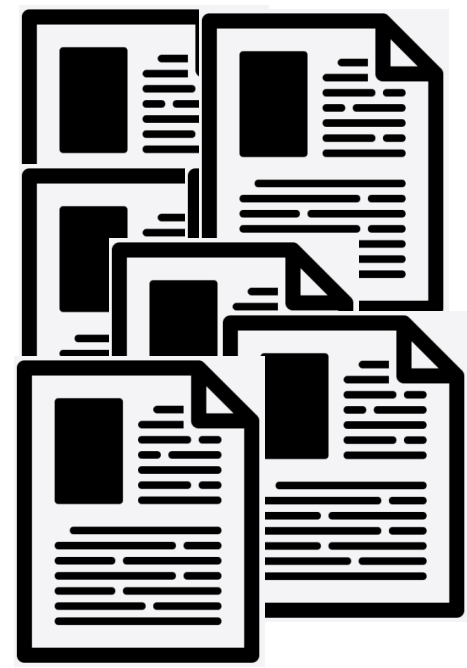


 Interact with User T



# Today: news recommendations

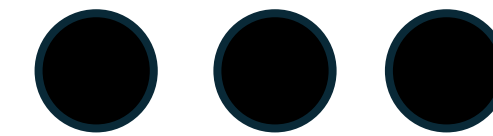
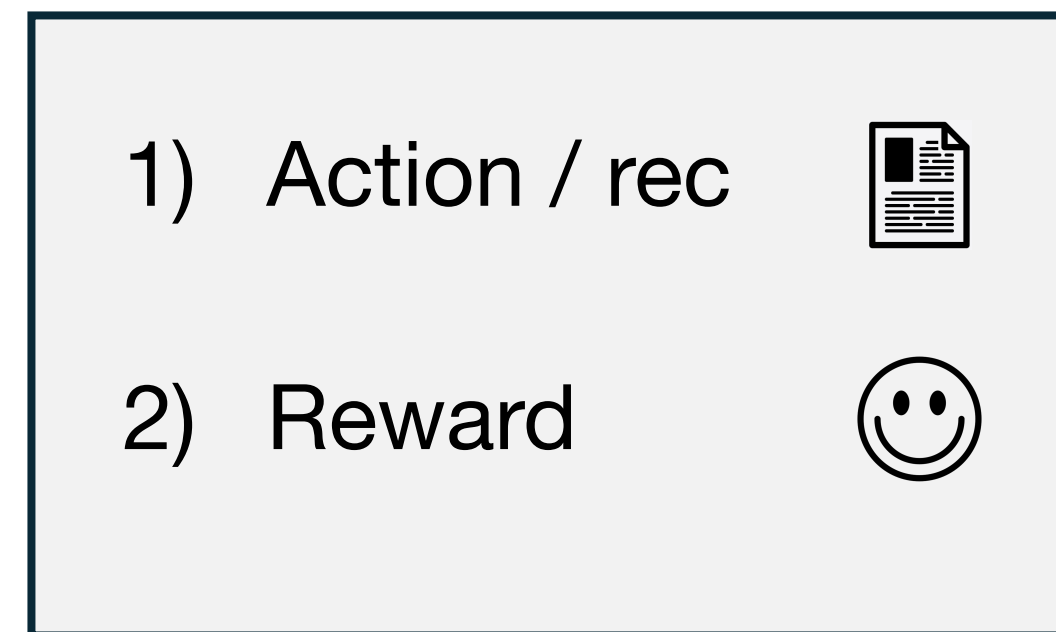
New articles are released



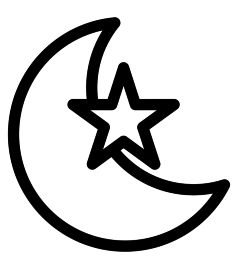
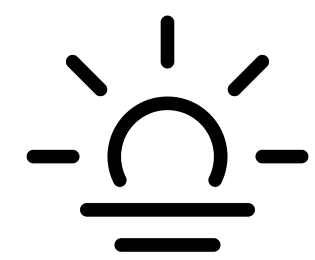
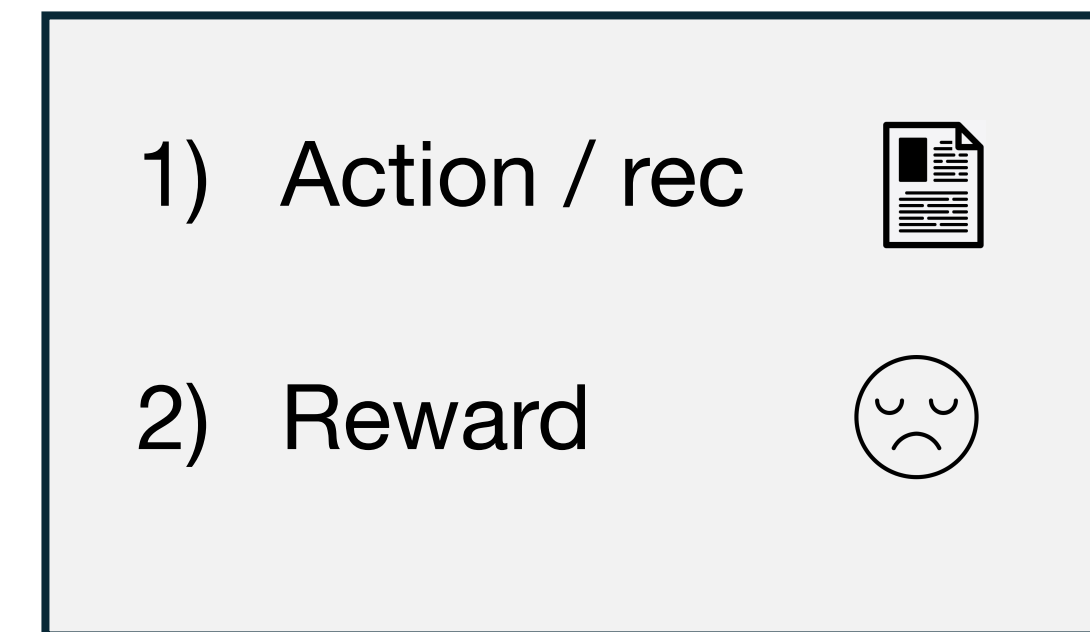
An LLM reads them



 Interact with User 1



 Interact with User T



Repeat process tomorrow

# Goal: sharpen beliefs with more info

- When articles are released, AI system must
  1. Form informed prior based on article text
  2. Gather data to resolve remaining uncertainty
  3. Balance exploration/exploitation
- For this talk, we assume users are **exchangeable**
  - Algo generalizes to personalized settings with user features

# Thompson sampling



$Z$  : Article features

$U$  : Other latent factors that govern article popularity

- Draw  $U$  from the posterior given all data about the article
- Pick best article according to the drawn values

# Thompson sampling

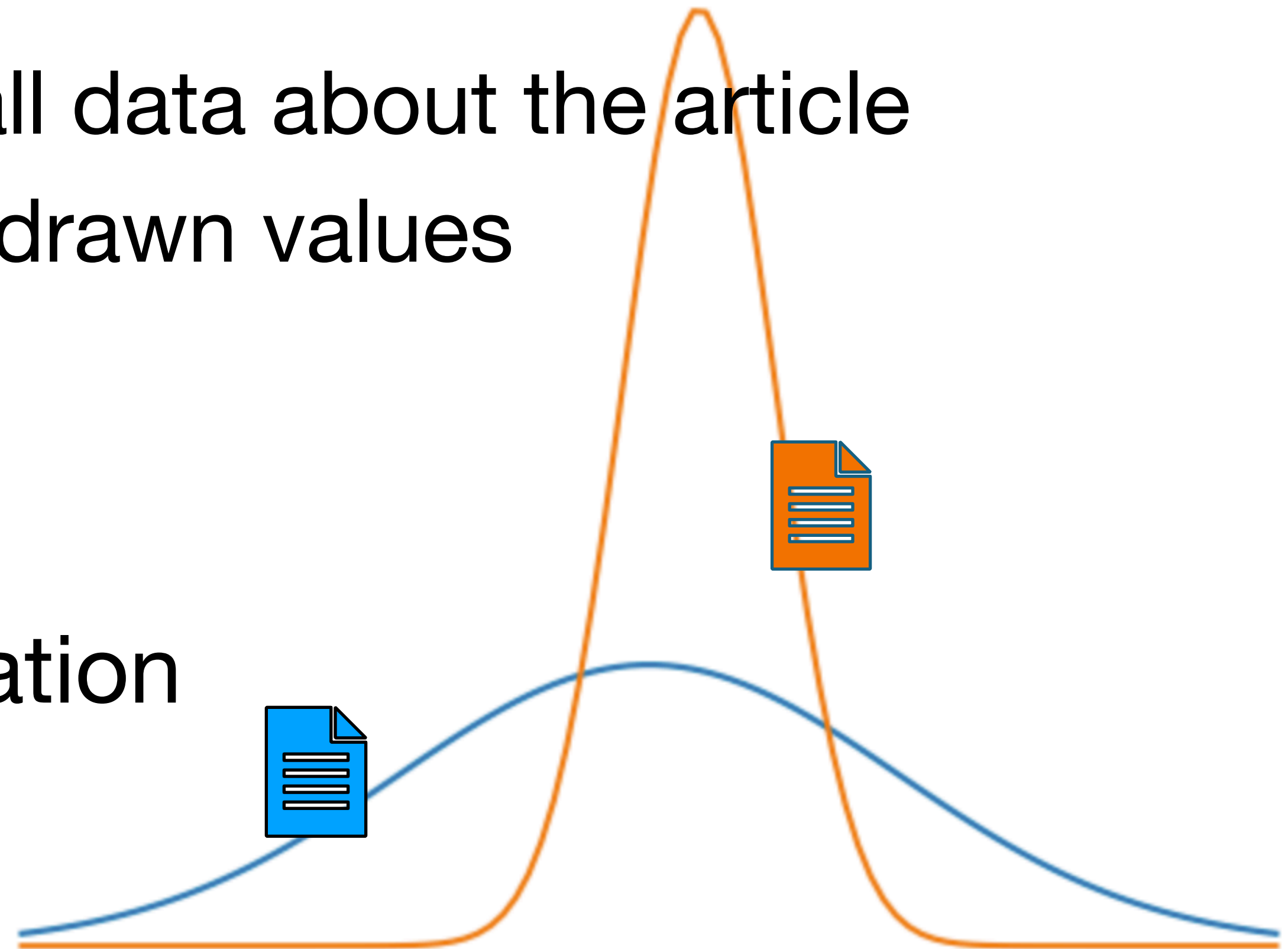


$Z$  : Article features

$U$  : Other latent factors that govern article popularity

- Draw  $U$  from the posterior given all data about the article
- Pick best article according to the drawn values

Balances exploration and exploitation



# Thompson sampling



$Z$  : Article features

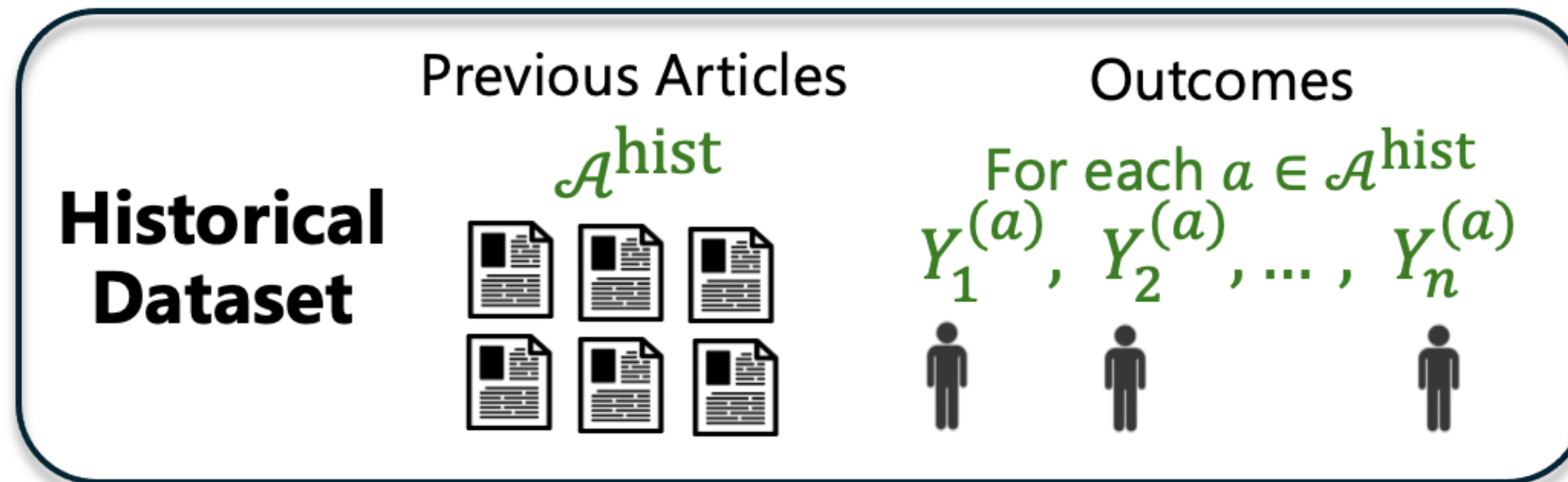
$U$  : Other latent factors that govern article popularity

- Draw  $U$  from the posterior given all data about the article
- Pick best article according to the drawn values

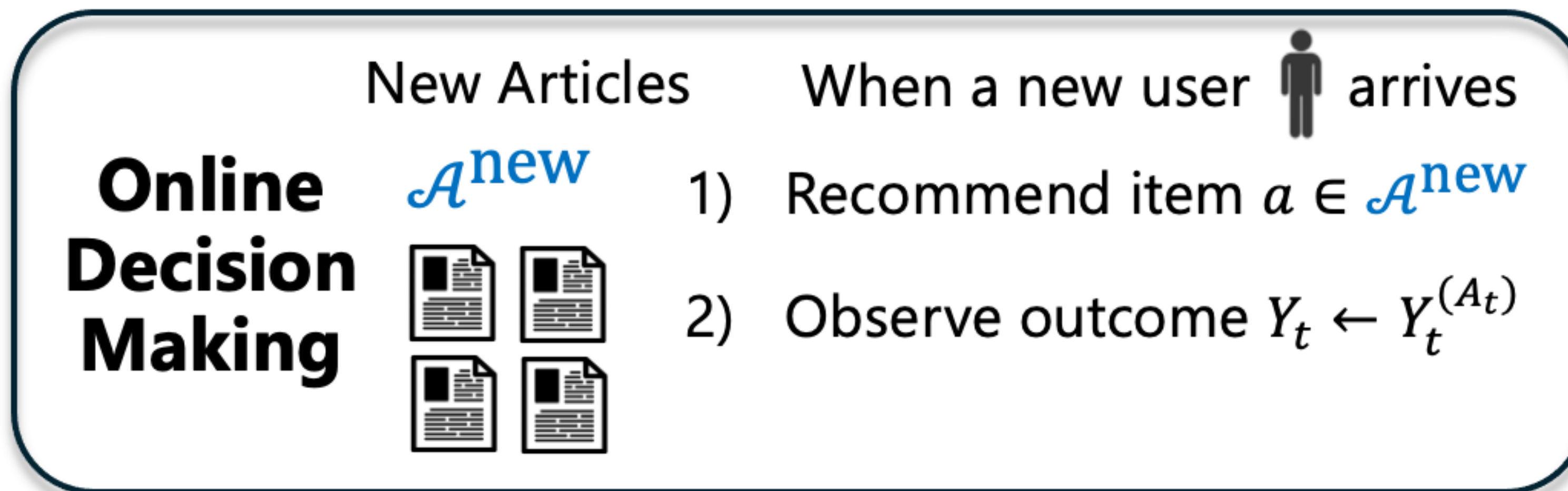
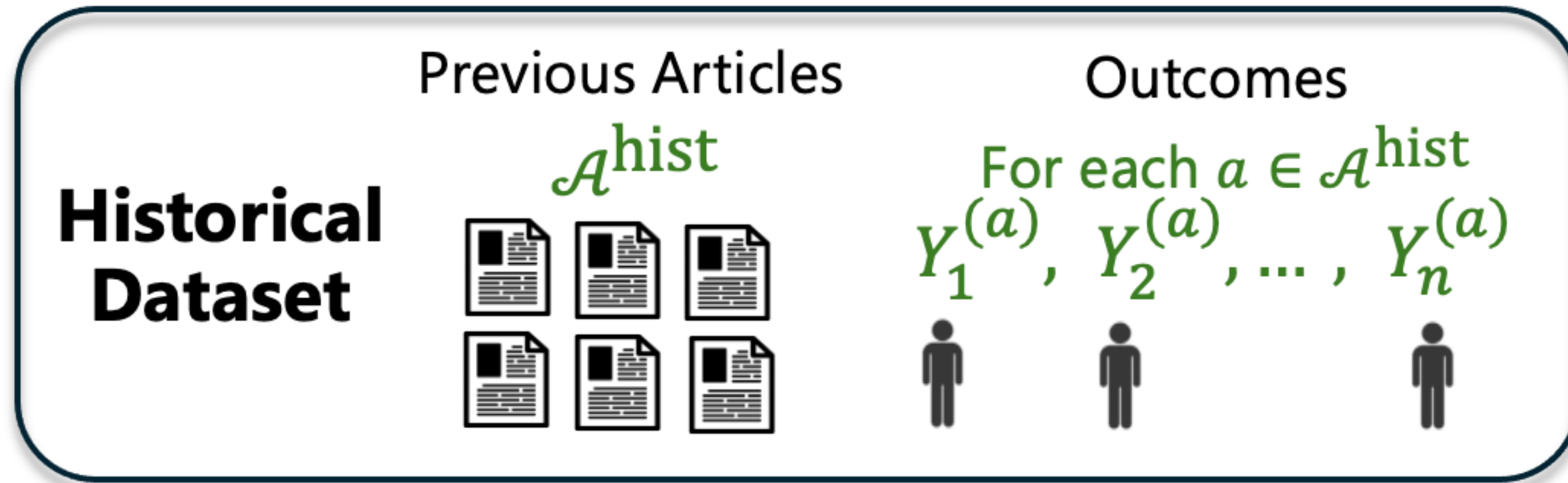
## Main challenge

Informed exploration requires probabilistic model over text

# Our solution: use massive historical data



# Step 1: Pretrain a sequence model

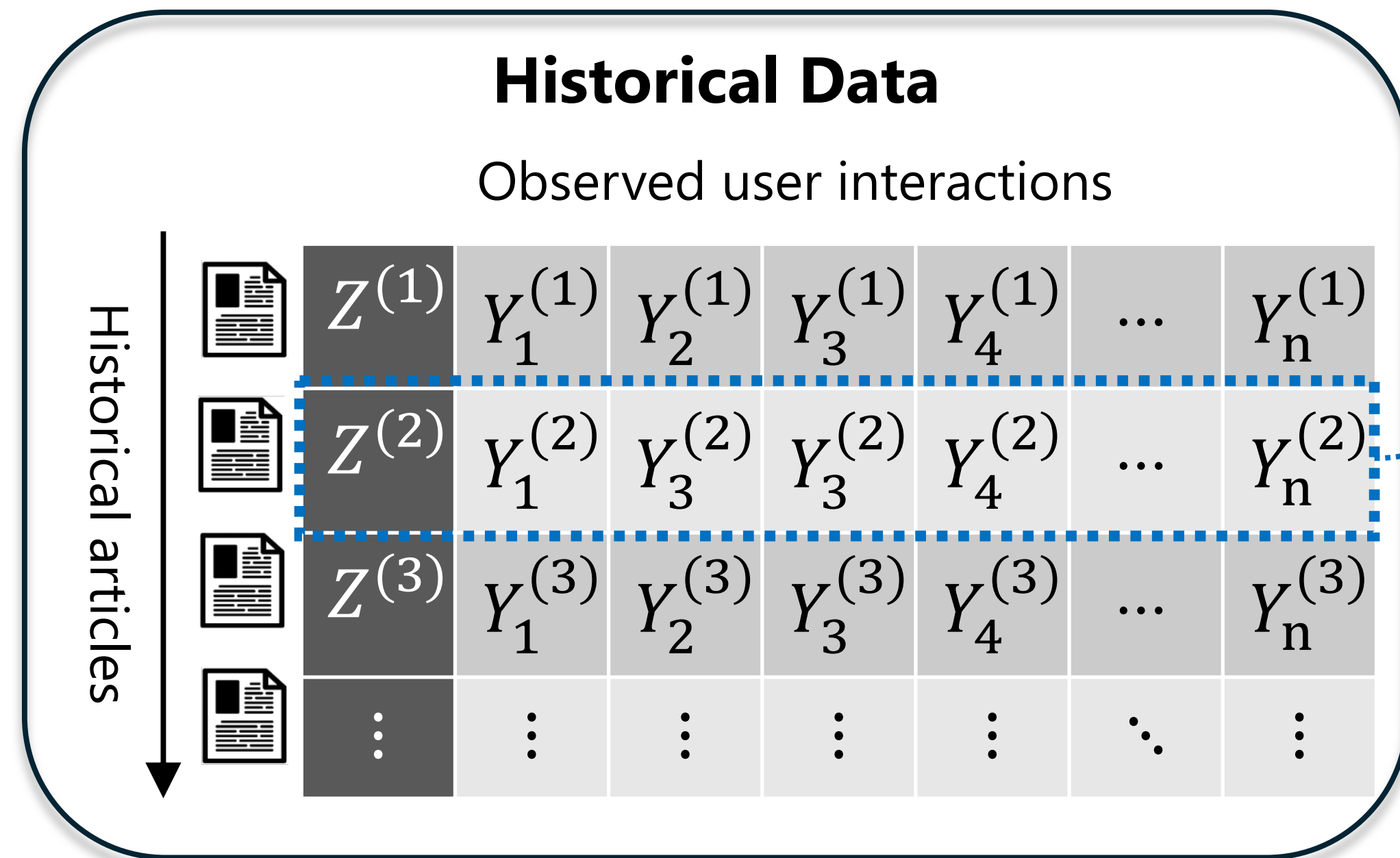


Use vast historical data to **warm-start** online decision-making

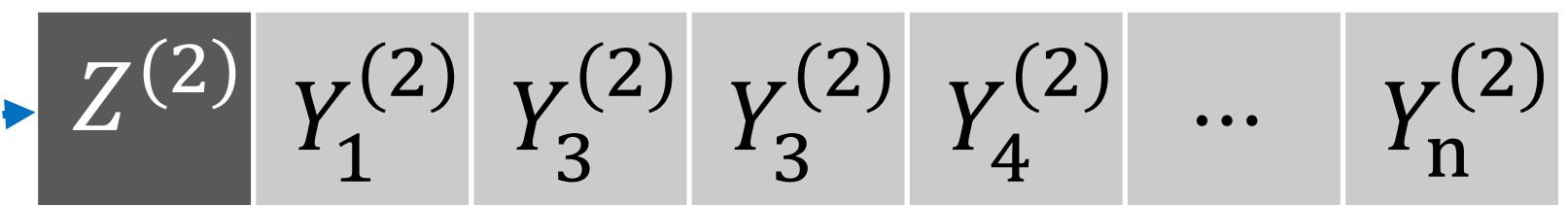


# Step 1: Pretrain a sequence model

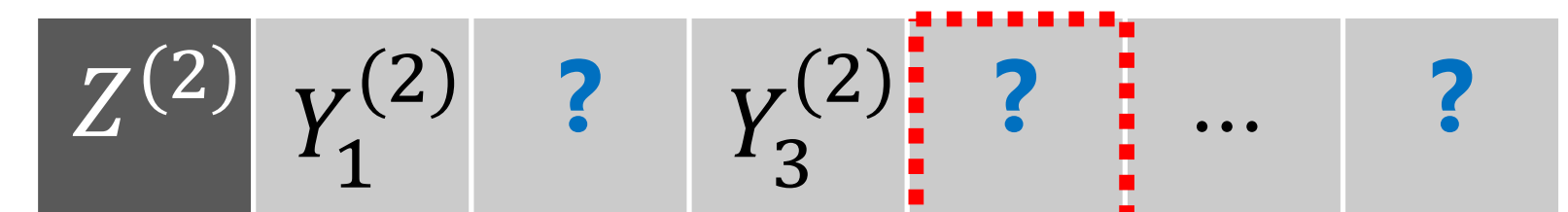
Low loss requires ability to sharpen beliefs



**A) Pick an article at random**



**B) Mask & predict some user interactions**

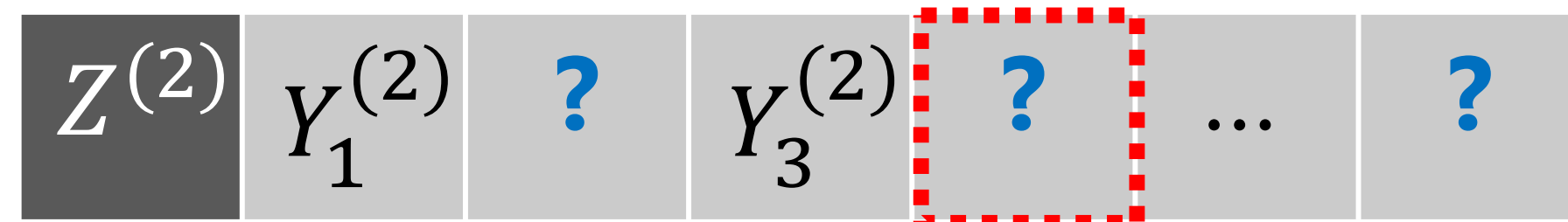


$$P(y \mid Z^{(2)}, Y_1^{(2)}, Y_3^{(2)}) = ?$$

# Step 1: Pretrain a sequence model

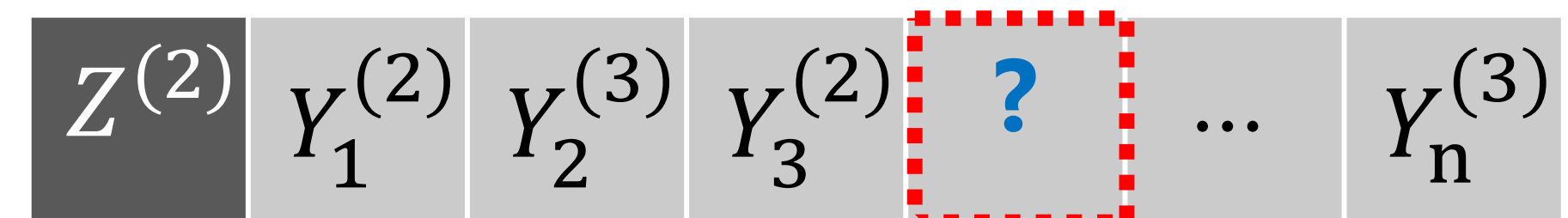
Low loss requires ability to sharpen beliefs

Predict well having observed few Y's



$$P(y \mid Z^{(2)}, Y_1^{(2)}, Y_3^{(2)}) = ?$$

Predict well having observed more Y's



$$P(y \mid Z^{(2)}, Y_1^{(2)}, Y_2^{(2)}, Y_3^{(2)}, Y_n^{(2)}) = ?$$

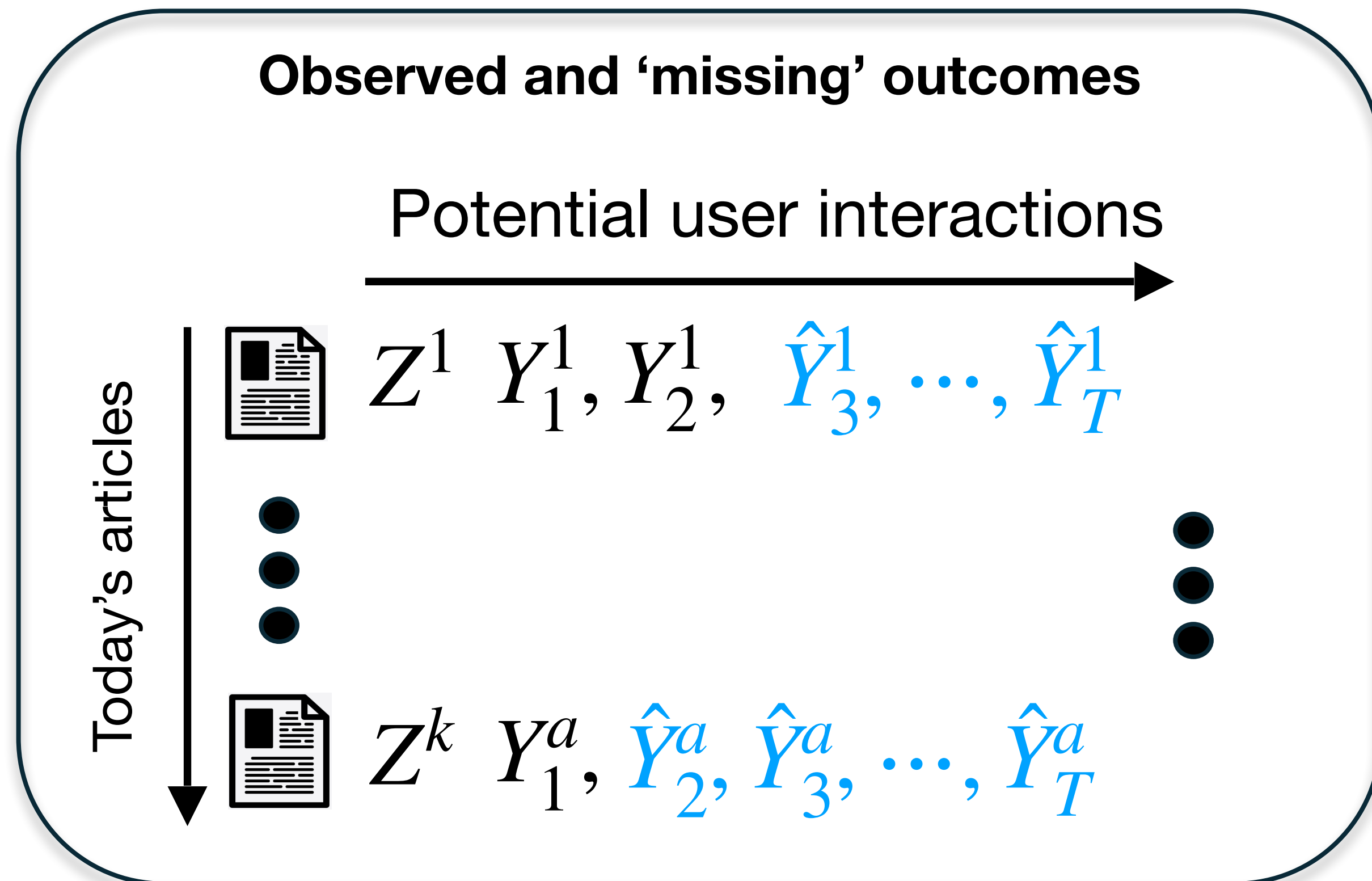
**Sequence predictions enable probabilistic reasoning**

1. When to heavily weigh the “prior” based on text Z?
2. When user interactions overwhelm the “prior”?

# Step 2: Act to resolve uncertainty

Autoregressive generation reveals actions that *might* have great performance

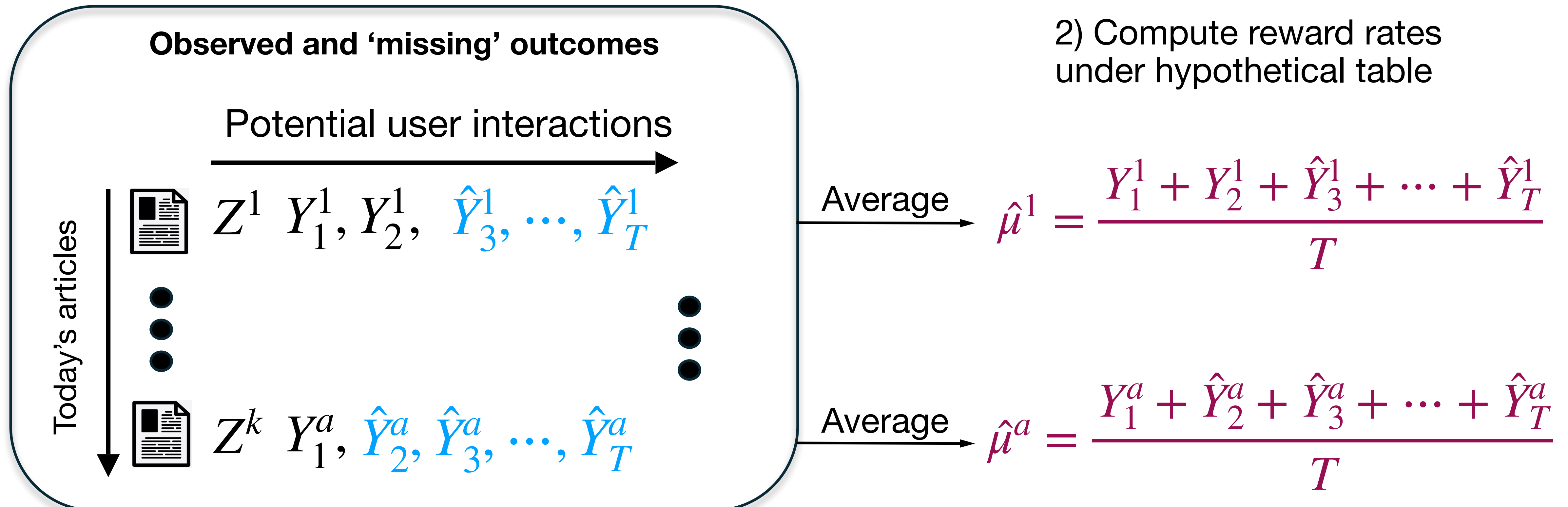
- 1) Fill in missing outcomes by autoregressive generation



# Step 2: Act to resolve uncertainty

Autoregressive generation reveals actions that *might* have great performance

1) Fill in missing outcomes by autoregressive generation



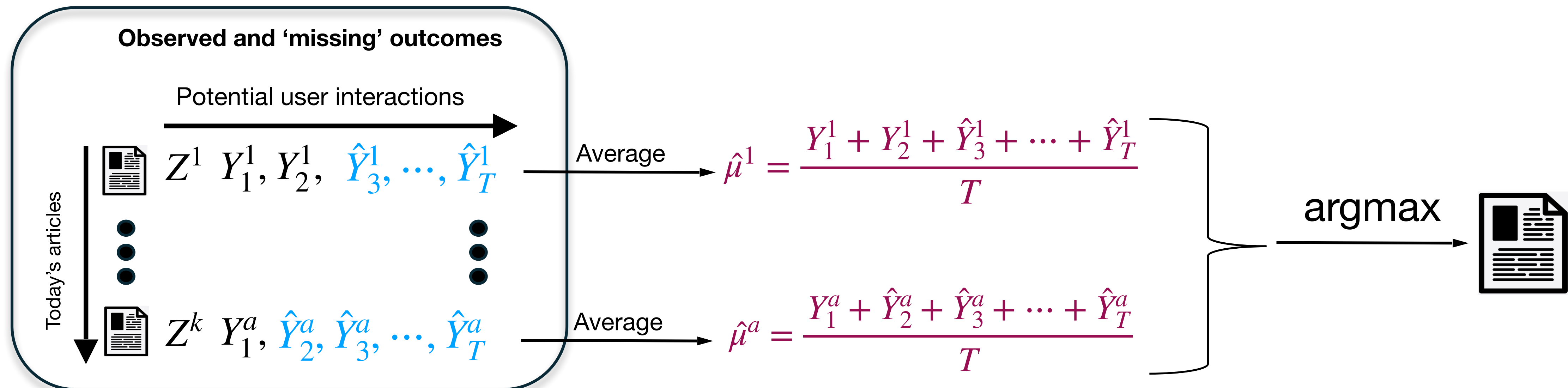
# Step 2: Act to resolve uncertainty

Autoregressive generation reveals actions that *might* have great performance

1) Fill in missing outcomes by autoregressive generation

2) Compute reward rates under hypothetical table

3) Pick item with highest hypothetical reward



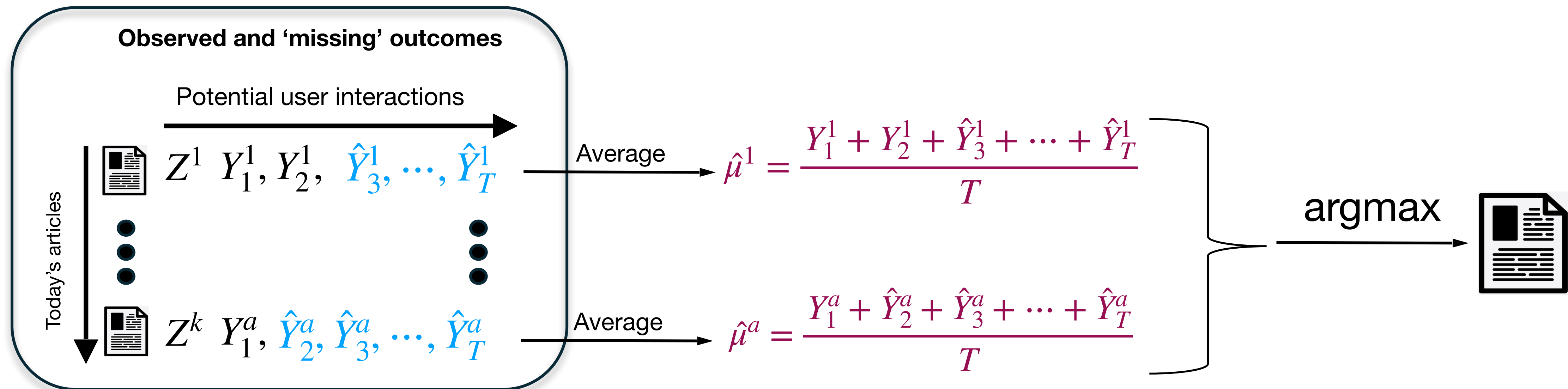
# Step 2: Act to resolve uncertainty

Autoregressive generation reveals actions that *might* have great performance

1) Fill in missing outcomes by autoregressive generation

2) Compute reward rates under hypothetical table

3) Pick item with highest hypothetical reward



**Action with high potential has a fair chance**

# Theorem: regret controlled by perplexity

## Data assumptions

### 1. Independently drawn articles.

- Text/ potential outcome sets  $(Z^{(a)}, Y_1^{(a)}, \dots, Y_T^{(a)})$  are independent across articles.

### 2. Historical data is representative of future days

- Distribution  $(Z^{(a)}, Y_1^{(a)}, \dots, Y_T^{(a)})$  is the same for articles in historical data as what governs tomorrow's draw.

### 3. Exchangeability across users

- $$P^* \left( Y_1^{(a)}, \dots, Y_T^{(a)} \mid Z^{(a)} \right) = P^* \left( Y_{\sigma(1)}^{(a)}, \dots, Y_{\sigma(T)}^{(a)} \mid Z^{(a)} \right)$$

# Theorem: regret controlled by perplexity

Optimal decision under imagined data = best under posterior draw

Two assumptions: 1) bounded rewards, 2) training length sequence exceeds T

Regret when using  
autoregressive model  $p_\theta$

$\leq$

Regret of Thompson  
sampling with true prior

$$+ \sqrt{2 \cdot \text{no. articles each day} \cdot (\ell_T(p^*) - \ell_T(p_\theta))}$$

Optimal autoregressive  
prediction

Validation  
loss



# Theorem: regret controlled by perplexity

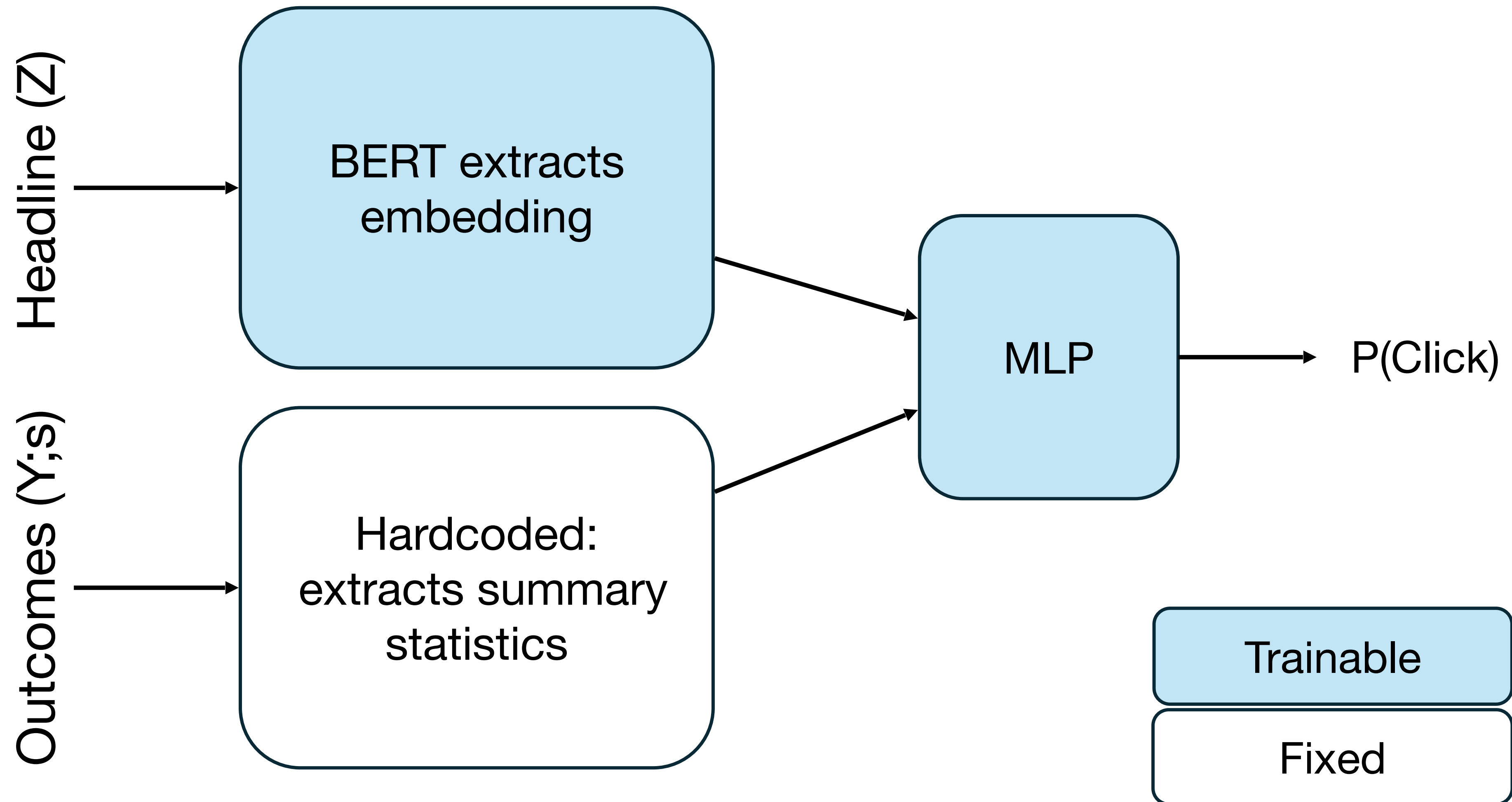
Optimal decision under imagined data = best under posterior draw

Two assumptions: 1) bounded rewards, 2) training length sequence exceeds  $T$

$$\begin{aligned} \text{Regret when using} \\ \text{autoregressive model } p_\theta &\leq \text{Regret of Thompson} \\ &\text{sampling with true prior} \\ &+ \sqrt{2 \cdot \text{no. articles each day} \cdot (\ell_T(p^\star) - \ell_T(p_\theta))} \end{aligned}$$

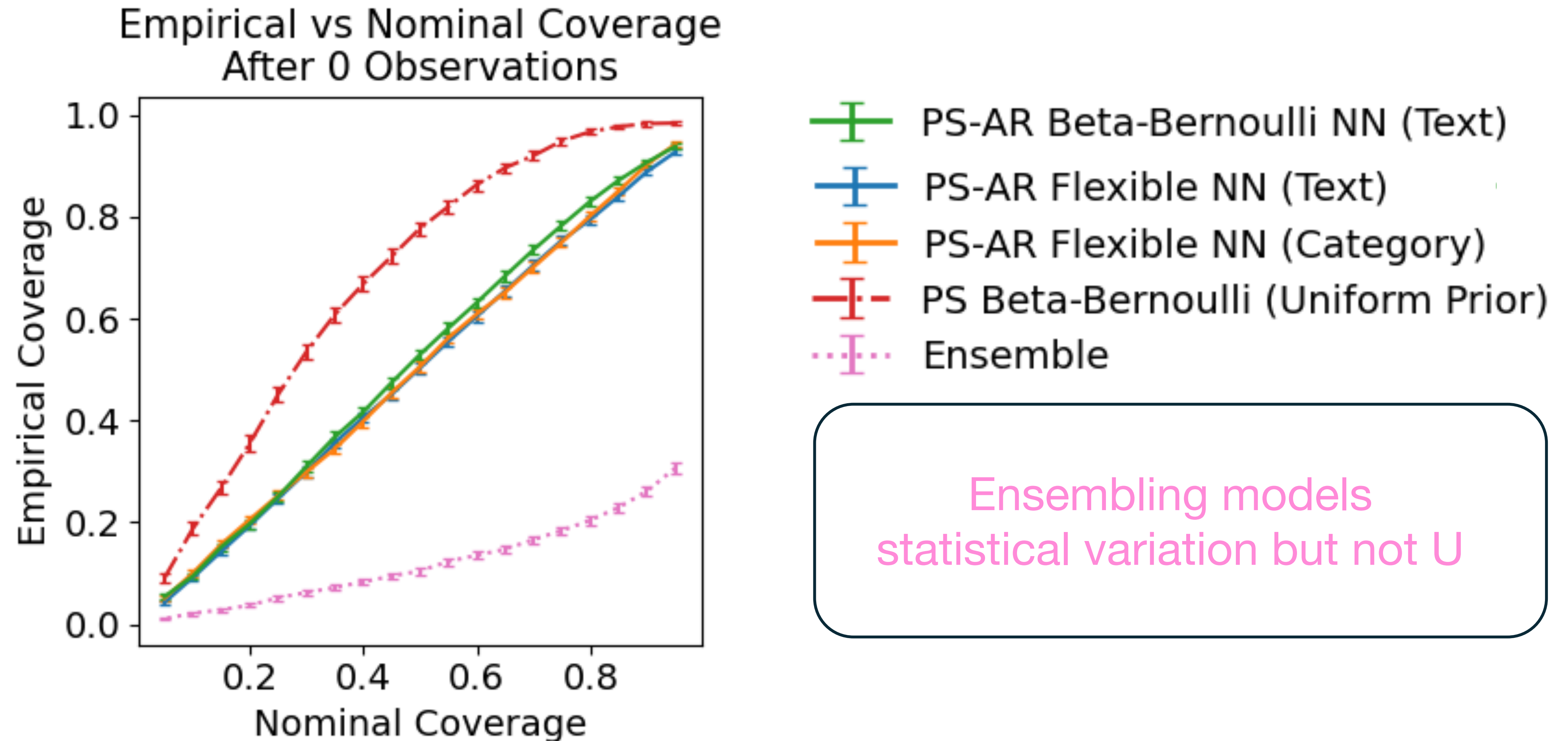
Scaling laws over user interactions control recommendation quality!

# Proof of concept: baby sequence model



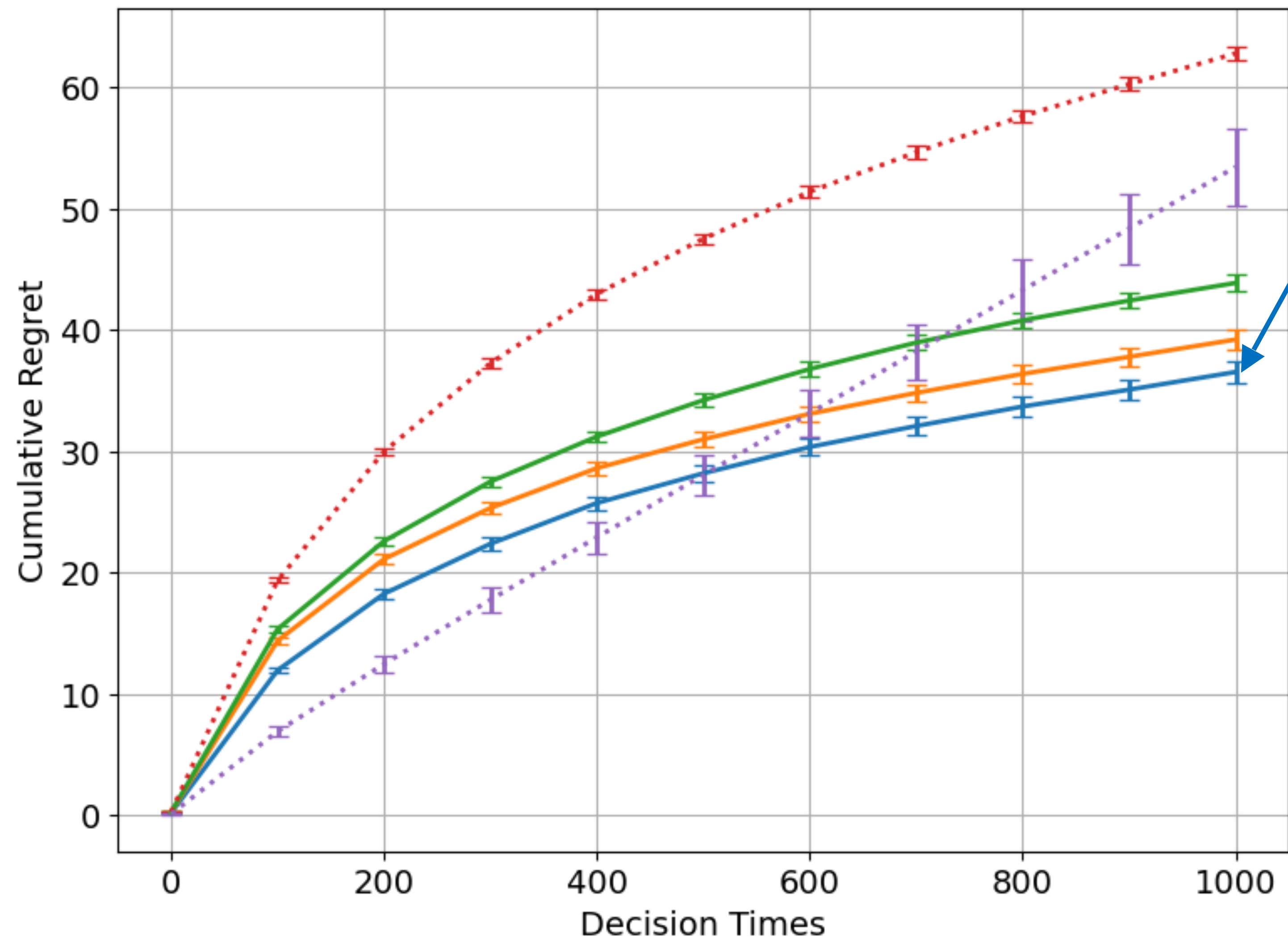
# Coverage

Autoregressive generation mimics proper Bayesian beliefs given headline (text)



# Regret

A semi-realistic simulator using public MSN news article data



Autoregressive generation scales to settings that requires end-to-end finetuning an LLM.

- PS AutoReg Flexible Neural Net, Text
- PS AutoReg Flexible Neural Net, Category
- PS Beta-Bernoulli (Uniform Prior)
- UCB
- Greedy

# **Vision: probabilistic reasoning**

**Intelligent agents must comprehend uncertainty and take actions to resolve it**

**Going beyond web-data requires  
careful data collection**

1. Formulate informed prior
2. Decide which data to collect  
subject to cost constraints
3. Update beliefs

# Vision: probabilistic reasoning

Intelligent agents must comprehend uncertainty and take actions to resolve it

**Going beyond web-data requires careful data collection**

1. Formulate informed prior
2. Decide which data to collect subject to cost constraints
3. Update beliefs

**Proposal: Model-based planning through sequence models**

1. Rollout: autoregressively generate hypothetical data
2. Simulate posterior update by appending data to context
3. Guide real decision based on how sharply belief shifts

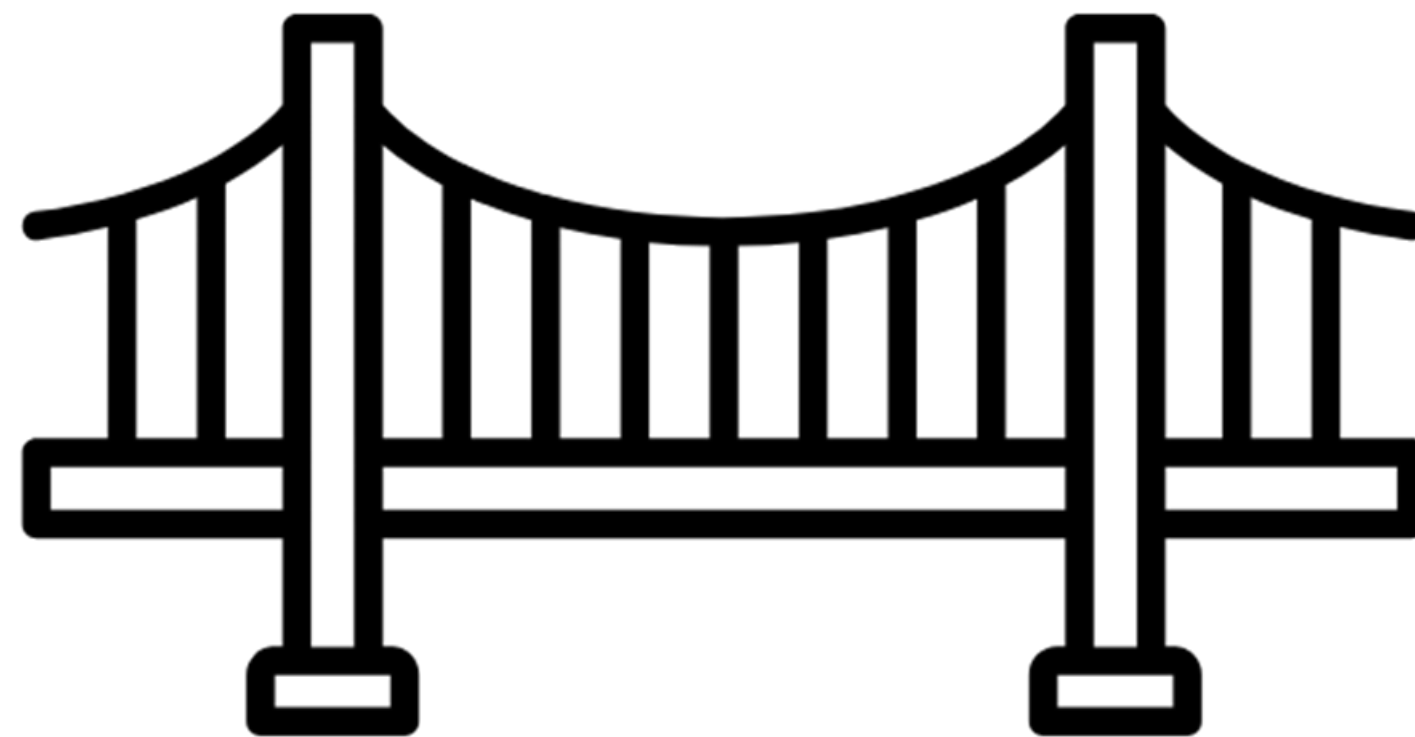
**Conceptual:** a well motivated problem crystalizing the insights

**Algorithmic:** link with interactive decision-making

**Theory:** accurate sequence modeling implies low regret

**Experiments:** scalable implementations with LLMs.

Generative  
Sequence  
Modeling



Exploration &  
Uncertainty  
Quantification

Posterior Sampling via Autoregressive Generation, ZCNR, [arXiv:2405.19466](https://arxiv.org/abs/2405.19466)

Exchangeable Sequence Models Quantify Uncertainty Over Latents, YN, [arXiv:2408.03307](https://arxiv.org/abs/2408.03307)