

In-Context Learning as Bayesian Inference

Tea Talk

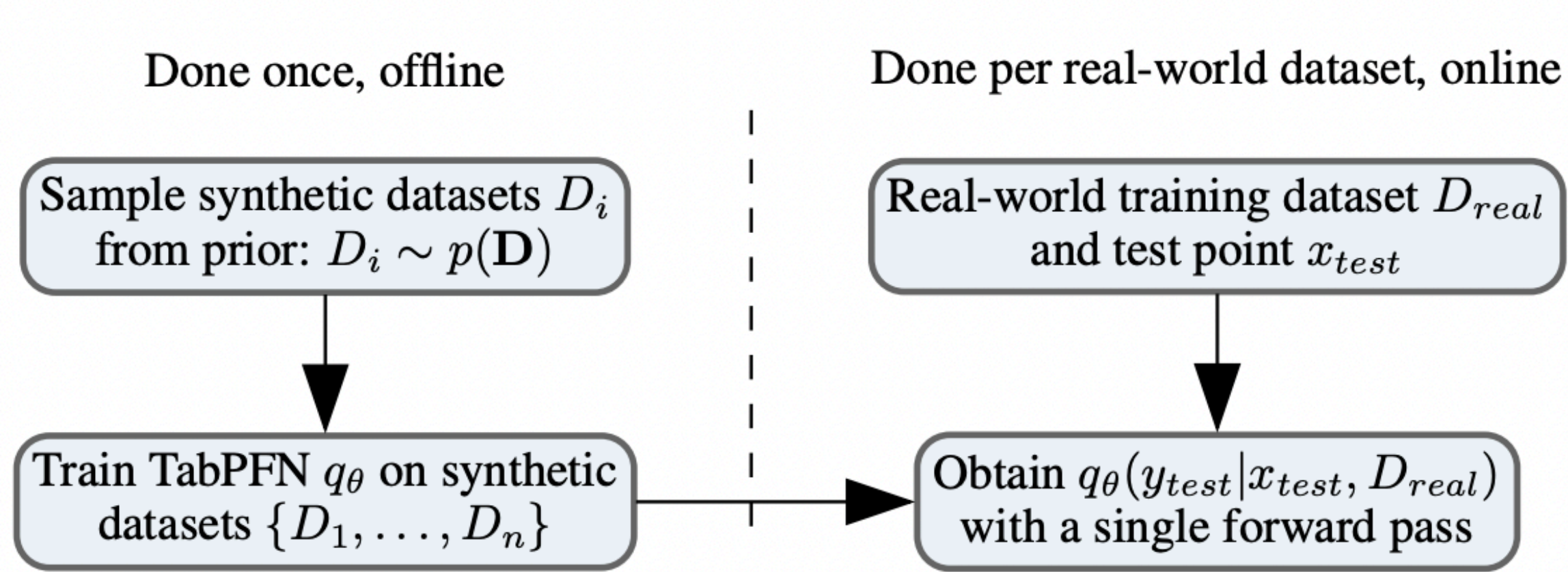
TABPFN: A TRANSFORMER THAT SOLVES SMALL TABULAR CLASSIFICATION PROBLEMS IN A SECOND

Noah Hollmann^{*,1,2} Samuel Müller^{*,1} Katharina Eggenberger¹ Frank Hutter^{1,3}

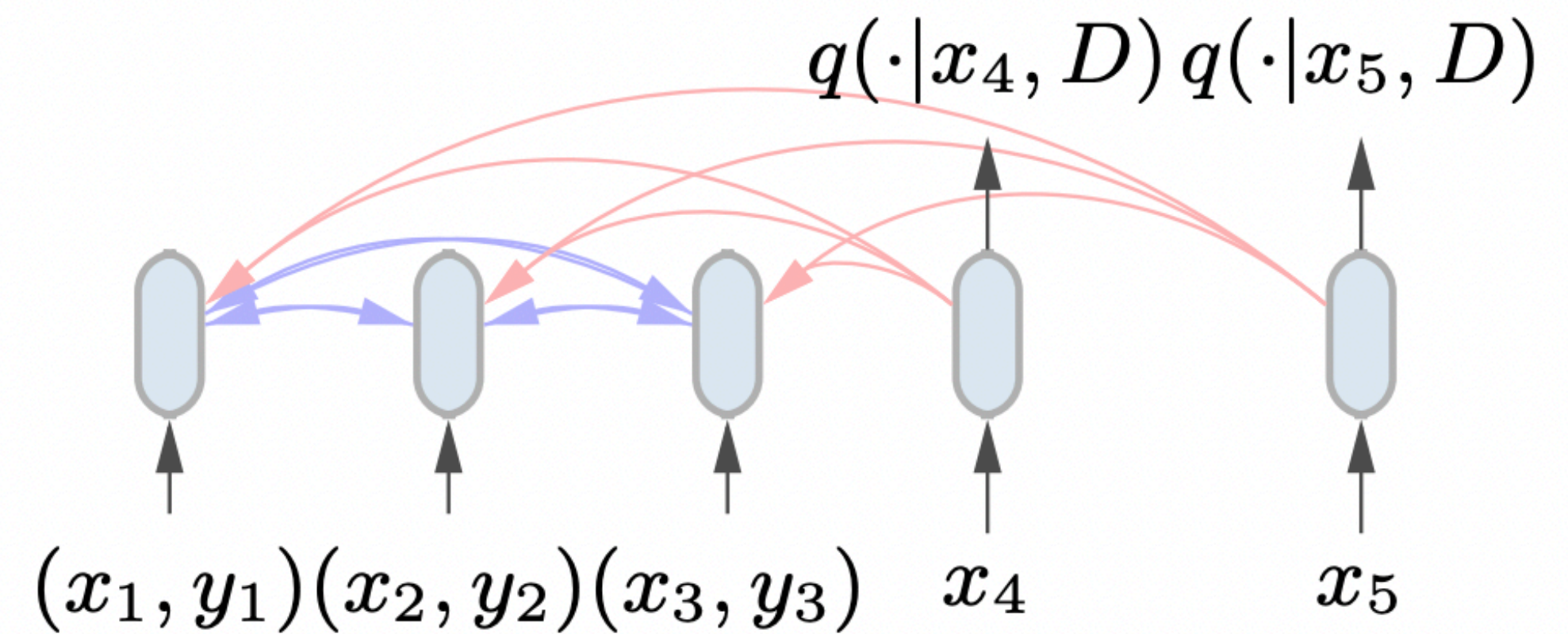
¹ University of Freiburg, ² Charité University Medicine Berlin

³ Bosch Center for Artificial Intelligence * Equal contribution.

Correspondance to noah.hollmann@charite.de & muellesa@cs.uni-freiburg.de



(a) Prior-fitting and inference



(b) Architecture and attention mechanism

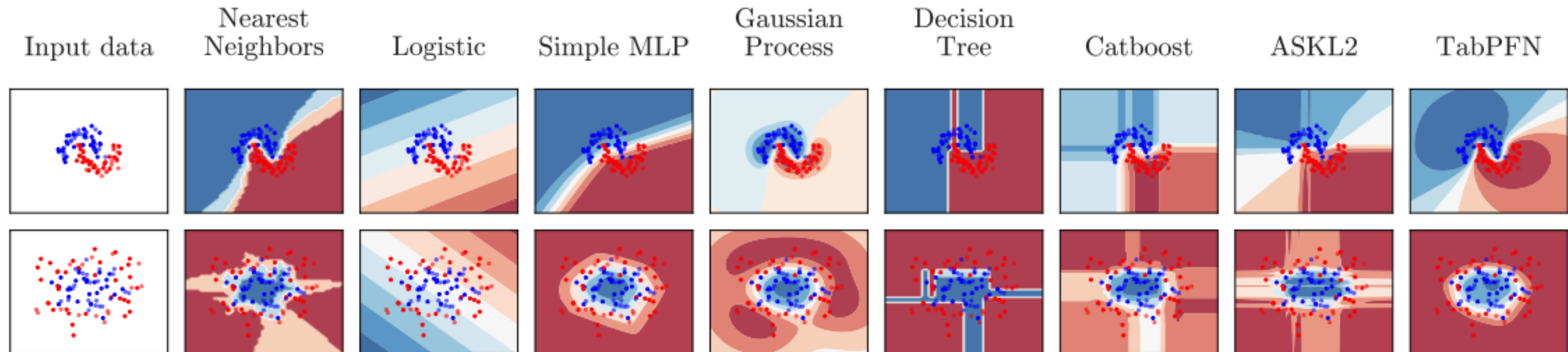
TABPFN: A TRANSFORMER THAT SOLVES SMALL TABULAR CLASSIFICATION PROBLEMS IN A SECOND

Noah Hollmann^{*,1,2} Samuel Müller^{*,1} Katharina Eggenberger¹ Frank Hutter^{1,3}

¹ University of Freiburg, ² Charité University Medicine Berlin

³ Bosch Center for Artificial Intelligence * Equal contribution.

Correspondance to noah.hollmann@charite.de & muellesa@cs.uni-freiburg.de



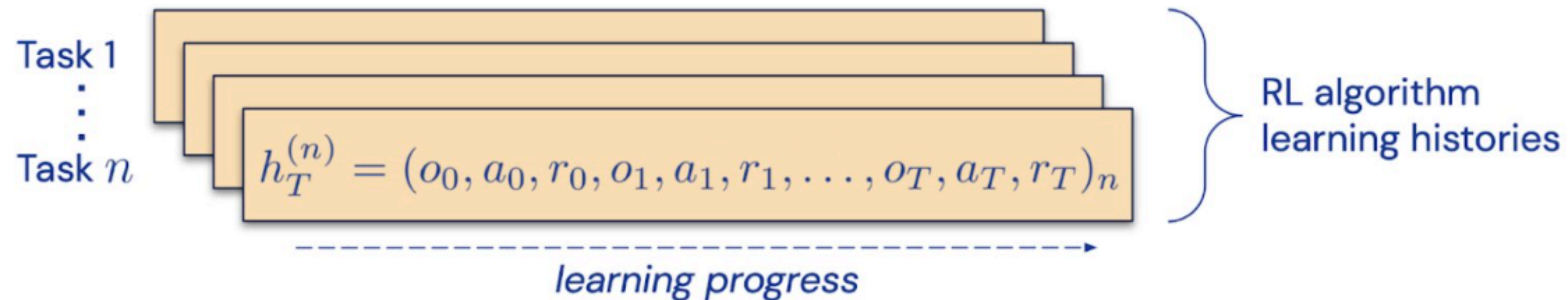
IN-CONTEXT REINFORCEMENT LEARNING WITH ALGORITHM DISTILLATION

DeepMind

First, we collect a dataset of learning histories from an RL algorithm trained on diverse tasks.

This can be *any* RL algorithm – it can be doing gradient updates, replay, planning, can be on or off-policy, model-free or model-based.

Data Generation



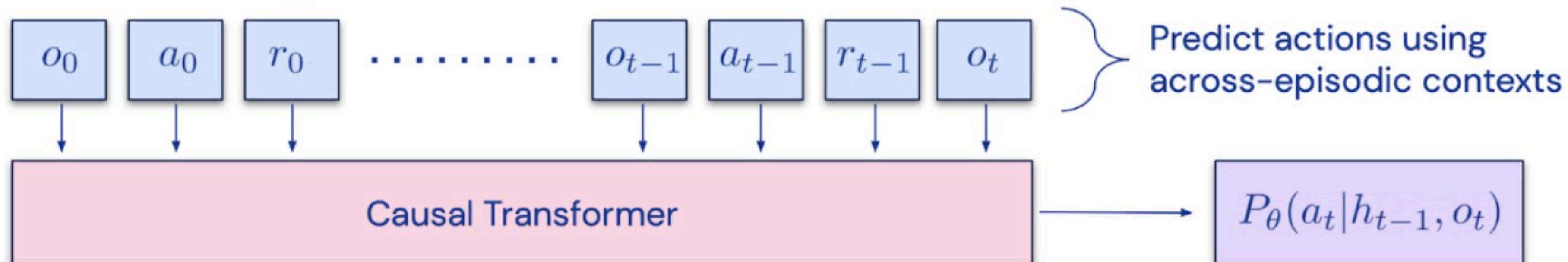
IN-CONTEXT REINFORCEMENT LEARNING WITH ALGORITHM DISTILLATION

DeepMind

Next, train a transformer to predict actions from the entire learning history preceding the current timestep.

Policy *improves* throughout RL training, to predict actions accurately, transformer needs to

Model Training



An Explanation of In-context Learning as Implicit Bayesian Inference

Sang Michael Xie
Stanford University
xie@cs.stanford.edu

Aditi Raghunathan
Stanford University
aditir@stanford.edu

Percy Liang
Stanford University
pliang@cs.stanford.edu

Tengyu Ma
Stanford University
tengyuma@cs.stanford.edu

Traditional fine-tuning (not used for GPT-3)

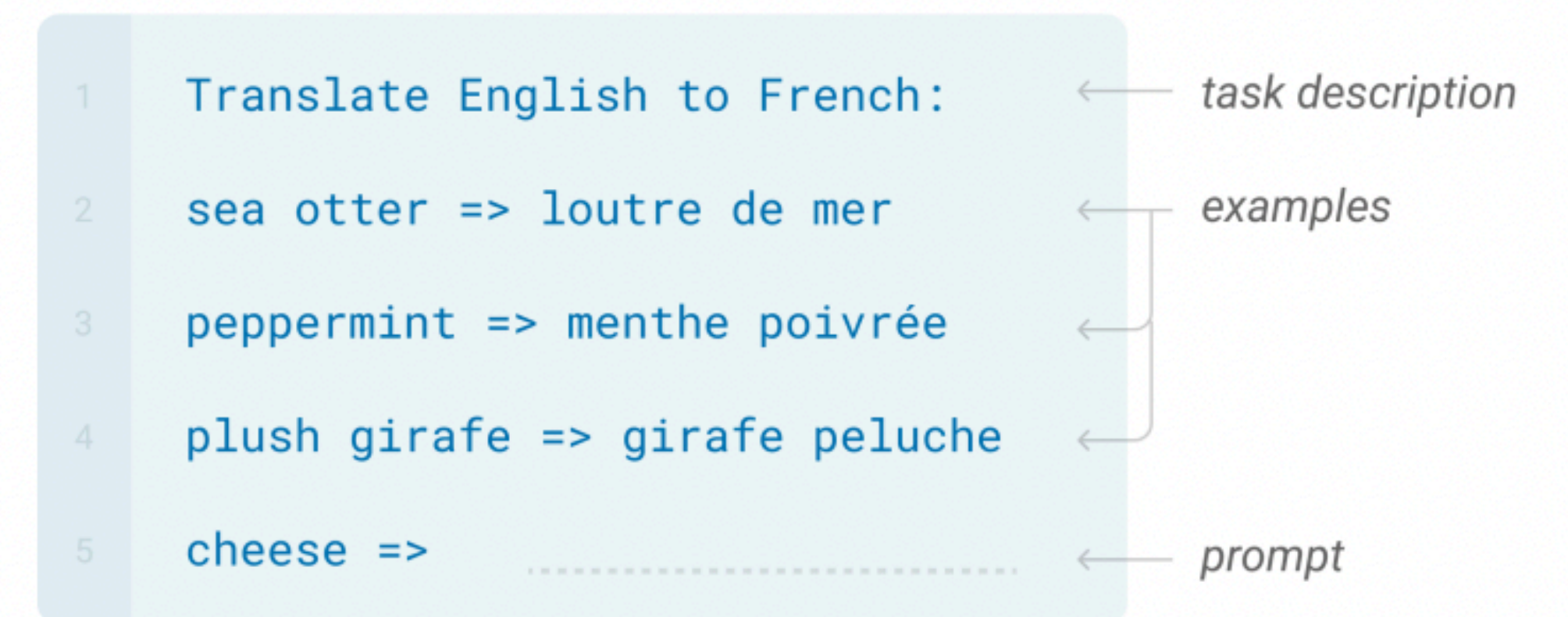
Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.

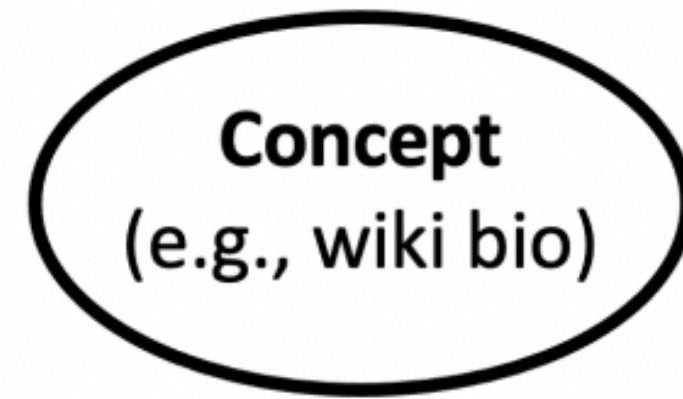


Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

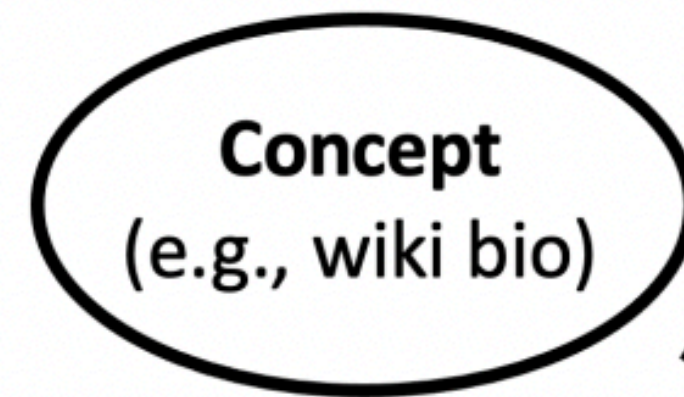


1. Pretraining documents are conditioned on a **latent concept** (e.g., biographical text)



Albert Einstein was a German theoretical physicist, widely acknowledged to be one of the greatest physicists of all time. Einstein is best known for developing the theory of relativity, but he also

2. Create independent examples from a **shared concept**. If we focus on full names, wiki bios tend to relate them to nationalities.



Input (x)	Output (y)	Delimiter
Albert Einstein was	German	\n
Mahatma Gandhi was	Indian	\n
Marie Curie was	?	...brilliant? ...Polish?

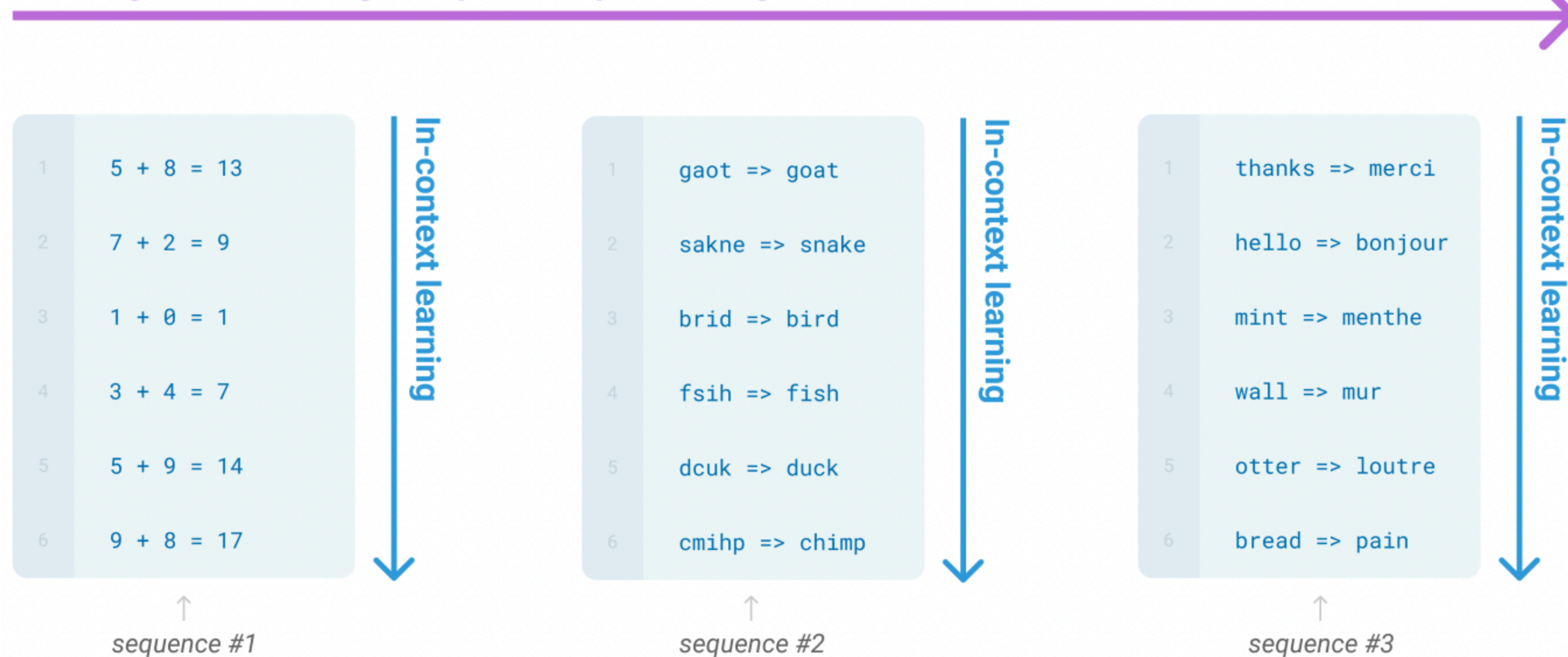
3. Concatenate examples into a prompt and predict next word(s). **Language model (LM)** implicitly infers the **shared concept** across examples despite the unnatural concatenation



Pretraining distribution. In our framework, a latent concept θ from a family of concepts Θ defines a distribution over observed tokens o from a vocabulary \mathcal{O} . To generate a document, we first sample a concept from a prior $p(\theta)$ and then sample the document given the concept. Each pretraining document is a length T sequence:

$$p(o_1, \dots, o_T) = \int_{\theta \in \Theta} p(o_1, \dots, o_T | \theta) p(\theta) d\theta. \quad (2)$$

Learning via SGD during unsupervised pre-training



$$p(\text{output}|\text{prompt}) = \int_{\text{concept}} p(\text{output}|\text{concept}, \text{prompt})p(\text{concept}|\text{prompt})d(\text{concept}).$$

Theorem 1. Assume the assumptions in Section 2.1 hold. If Condition 1 holds, then as $n \rightarrow \infty$ the prediction according to the pretraining distribution is

$$\arg \max_y p(y|S_n, x_{test}) \rightarrow \arg \max_y p_{prompt}(y|x_{test}). \quad (15)$$

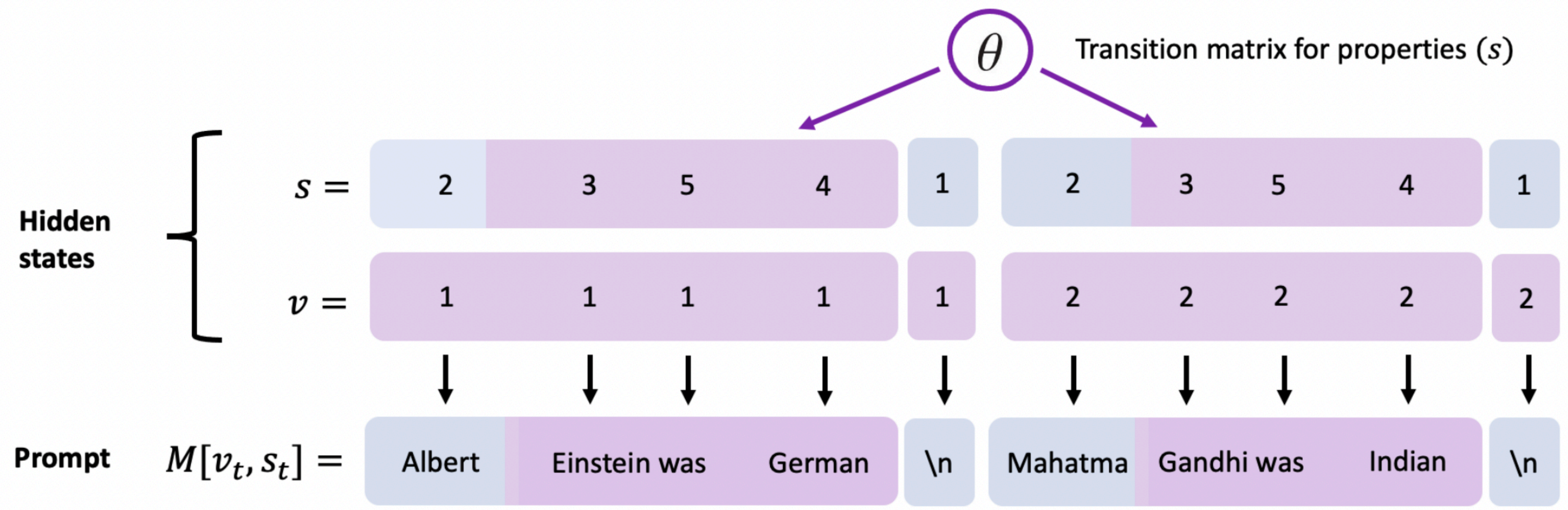
Thus, the in-context predictor f_n achieves the optimal 0-1 risk: $\lim_{n \rightarrow \infty} L_{0-1}(f_n) = \inf_f L_{0-1}(f)$.

$$[S_n, x_{test}] = [x_1, y_1, o^{\text{delim}}, x_2, y_2, o^{\text{delim}}, \dots, x_n, y_n, o^{\text{delim}}, x_{test}] \sim p_{prompt}.$$

Memory matrix $M =$

Entities (v)	Properties (s)					etc.
	Newline	First name	Last name	Nationality	Linking verb	
\n	Albert	Einstein	German	was		...
\n	Mahatma	Gandhi	Indian	was		
⋮						

Empirical HMM model (Prompt distribution)



Result

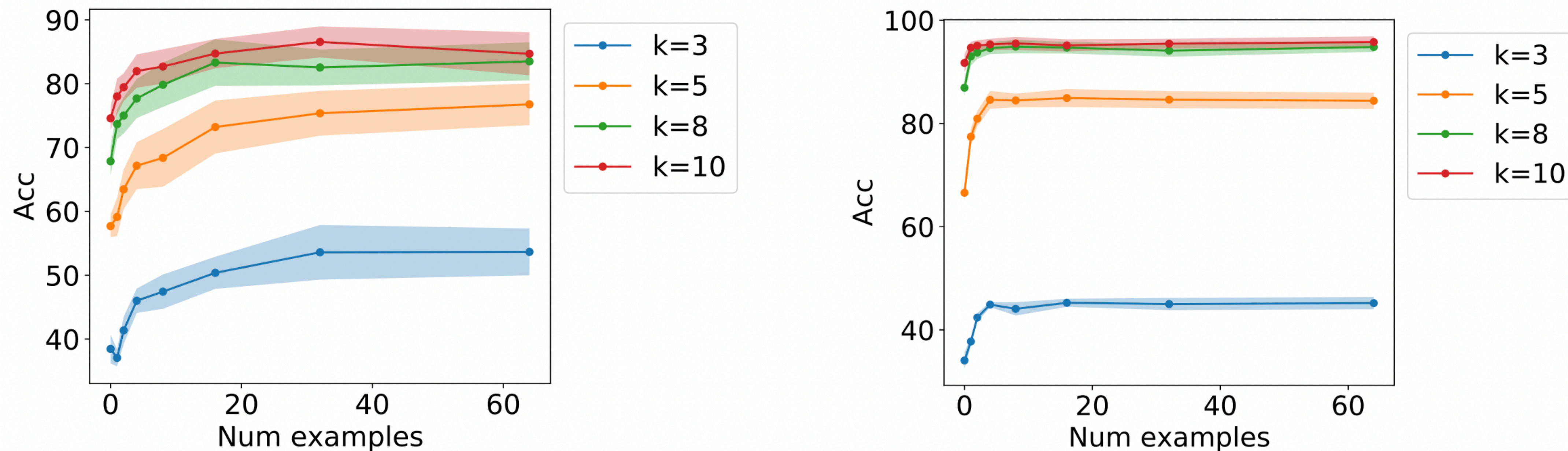


Figure 3: In-context accuracy (95% intervals) of Transformers (left) and LSTMs (right) on the GINC dataset. Accuracy increases with number of examples n and length of each example k .

Thank you

Questions

- No free lunch?