

Summary

Learn an **entailment-aware selective generator** with an abstaining option that controls the “**rate of hallucination**” with a probabilistic guarantee.

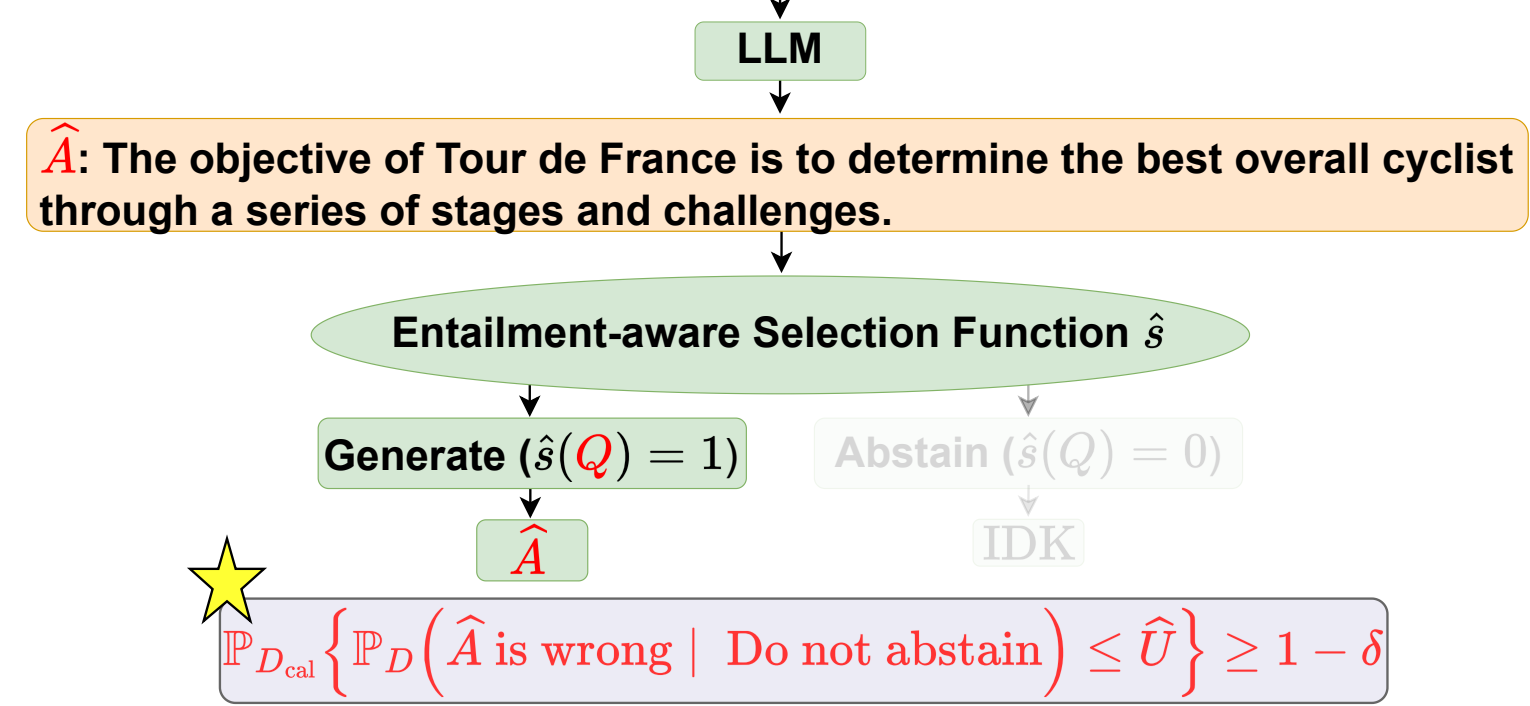
Problem: Hallucination Control in Language Generation

Goal: Learn a hallucination-controllable language generator.

Contributions

1. Propose the first “certified” selective generator for language models.
2. Leverage **textual entailment** as a **correctness metric**.
3. Design a **semi-supervised learning algorithm** for selective generation.
4. Prove a **controllability guarantee** of the proposed algorithm.

Q : What is the objective of Tour de France?
 A_{true} : The objective of Tour de France is bicycle race.)



Selective Classifier:

[Geifman & El-Yaniv, 2017]

$$\hat{S}(\mathbf{X}) := \begin{cases} \hat{\mathbf{Y}} & \hat{s}(\mathbf{X}) = 1, \\ \text{IDK} & \text{o.w.} \end{cases}$$

Selective Generator:

$$\hat{S}(Q) := \begin{cases} \hat{A} & \hat{s}(Q) = 1, \\ \text{IDK} & \text{o.w.} \end{cases}$$

Main Challenge: Metric Misalignment

Definition. Metric Misalignment

Learning Metric (e.g. EM) \neq Evaluation Metric (e.g. SC)

- Example:
 - Q : Where in the bible does it mention Sodom and Gomorrah?
 - A_{true} : The book of Genesis mentions Sodom and Gomorrah.
 - \hat{A} : The story of Sodom and Gomorrah is found in Genesis 19.
- A standard learning metric on correct answers, i.e. Exact Match (EM), assumes a **single** correct answer (i.e. $\hat{A} =_{\text{EM}} A_{\text{true}}?$)
- As $\hat{A} \neq_{\text{EM}} A_{\text{true}}$, \hat{A} is **wrong** even if it is **semantically correct** (SC) ☹.

Idea 1. Textual Entailment as a Correctness Metric

Definition. Correctness Metric by Entailment

A generated answer \hat{A} is correct if

$$\hat{A} \in E_{\text{true}}(A_{\text{true}}) := \{\tilde{A} \mid \tilde{A} \text{ entails } A_{\text{true}}\}.$$

Definition. False Discovery Rate with Entailment (FDR-E)

$$\text{Learning Metric: } \mathbb{P}_{\mathcal{D}} \left(\hat{A} \notin E_{\text{true}}(A_{\text{true}}) \mid \hat{S}(Q) \neq \text{IDK} \right)$$

- We find a learning algorithm to control the **FDR-E**.

Idea 2. Pseudo-labeling for Textual Entailment

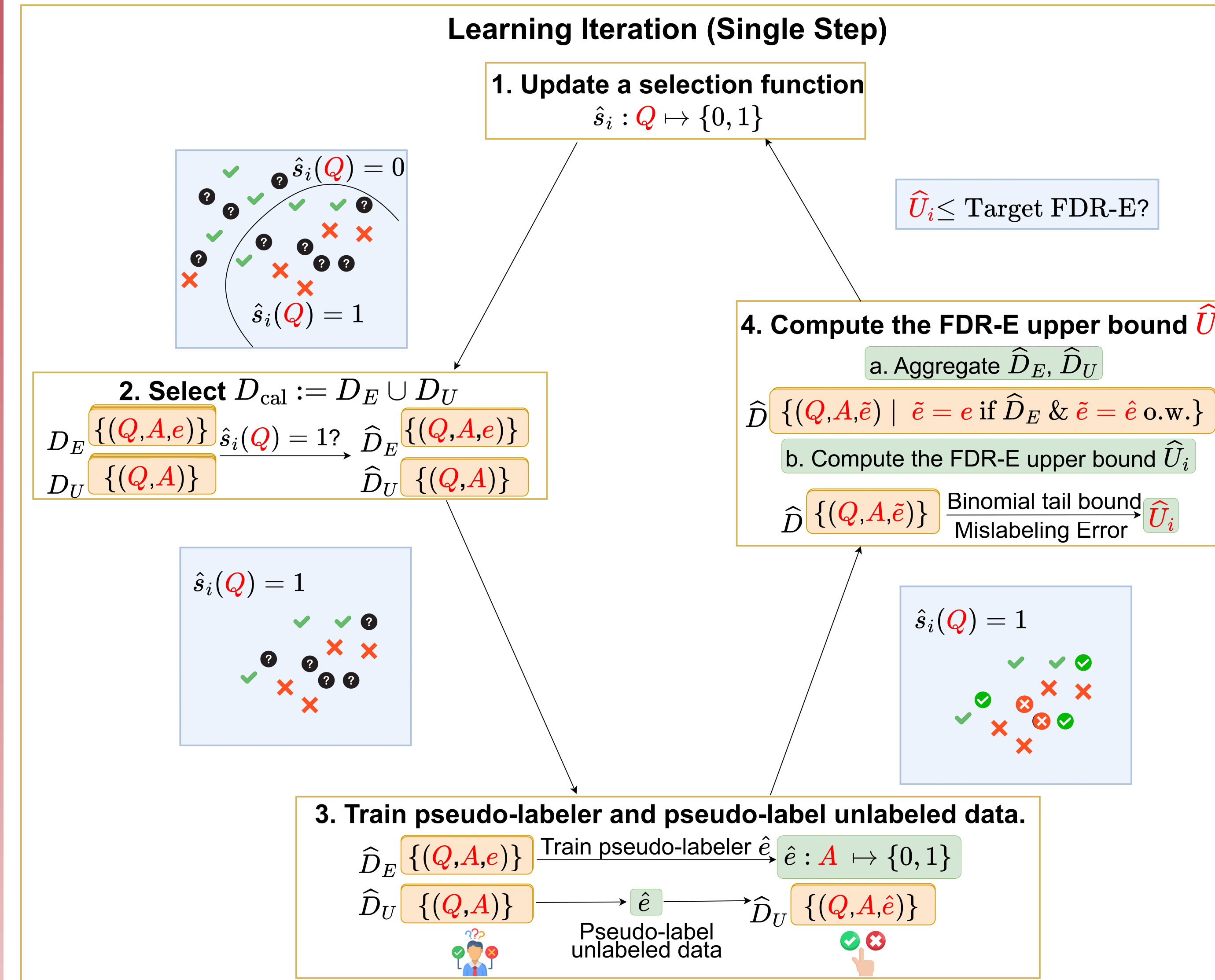
Calibration Set

$$\{(Q, A_{\text{true}}, \hat{A} \in E_{\text{true}}(A_{\text{true}}))\} \cup \{(Q, A_{\text{true}}, \hat{A} \in \hat{E}(A_{\text{true}}))\}$$

additional labels pseudo labels

- We propose a **label efficient semi-supervised learning algorithm**.

Solution: Semi-Supervised Selective Generator Learning

Algorithm SGen^{Semi} to find a selective generator \hat{S} Details on the Pseudo-Labeling Function \hat{E}

- We measure the estimation error by \hat{E} to upper-bound the FDR-E.

$$E_{\text{true}}(A_{\text{true}}) := \{\tilde{A} \mid \tilde{A} \text{ entails } A_{\text{true}}\} \text{ (True Entailment Set: Unknown)}$$

\hat{A}_1 : The book of Genesis Mentions Sodom and Gomorrah (A_{true}).
 \hat{A}_2 : The story of Sodom and Gomorrah can be found in Genesis 19 (\hat{A}).
 \hat{A}_3 : Genesis includes the story of Sodom and Gomorrah.
 \hat{A}_4 : Sodom and Gomorrah are referenced in the book of Genesis.

$\hat{E}(A_{\text{true}})$ (Estimated Entailment Set via Conformal Prediction)

$$\begin{aligned} \text{FDR-E (for SSL)} &= \mathbb{P}_{\mathcal{D}_{\hat{S}}} \{e = 0\} = \mathbb{P}_{\mathcal{D}_{\hat{S}}} \{e = 0, \hat{e} = 1\} \\ &- \mathbb{P}_{\mathcal{D}_{\hat{S}}} \{e = 1, \hat{e} = 0\} + \mathbb{P}_{\mathcal{D}_{\hat{S}}} \{\hat{e} = 0\} \end{aligned}$$

false neg.-entailment rate neg. entailment rate

- $e := \hat{A} \in E_{\text{true}}(A_{\text{true}})$
- $\hat{e} := \hat{A} \in \hat{E}(A_{\text{true}})$
- $\mathbb{P}_{\mathcal{D}_{\hat{S}}}(\cdot) := \mathbb{P}_{\mathcal{D}}(\cdot \mid \hat{S}(Q) \neq \text{IDK})$

Theoretical Result

Theorem. Controllability Guarantee on the FDR-E

For **any** LLMs and downstream language generation tasks, the following **model-agnostic** and **task-free controllability guarantee** holds:

$$\mathbb{P}_{\mathcal{D}_{\text{cal}}} \left\{ \underbrace{\mathbb{P}_{\mathcal{D}}(\hat{A} \notin E_{\text{true}}(A_{\text{true}}))}_{\hat{A} \text{ is "wrong"}} \mid \underbrace{\hat{S}(Q) \neq \text{IDK}}_{\text{Do not abstain}} \leq \hat{U} \right\} \geq 1 - \delta,$$

where δ is the confidence level and (\hat{s}, \hat{U}) is the algorithm output.

Experiments & Results

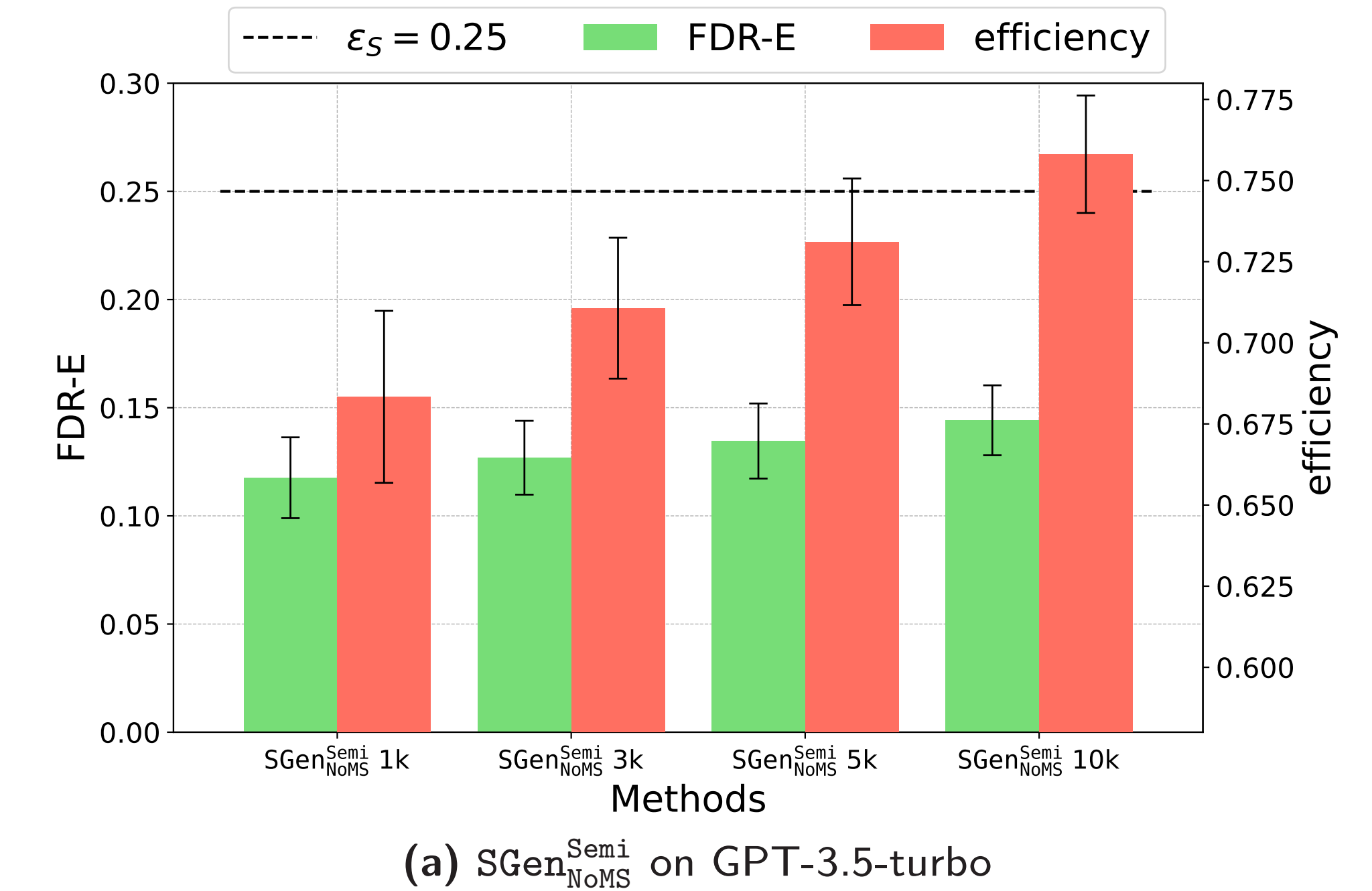
1. Benefit of Textual Entailment

- Our entailment-based learning metric shows better selection efficiency.
 - Selection efficiency: the proportion of non-abstained samples

Q	Who is the actor that plays Draco Malfoy?	When did the movie Benjamin Button come out?
A_{true}	Thomas Andrew Felton plays Draco Malfoy in the Harry Potter movies.	The movie Benjamin Button come out December 25, 2008.
\hat{A}	The actor who plays Draco Malfoy is Tom Felton. (correct)	The Curious Journey of Benjamin Button was released in 2008. (correct)
EM (Baseline)	rejected	rejected
Textual Entailment (Ours)	accepted	accepted

2. Benefit of Semi-Supervised Learning

- As we use more unlabeled data, selection efficiency gets better.



3. Benefit of Neuro-Selection Function

- Learning a selection function combination improves selection efficiency.

Models		GPT-3.5-turbo		Alpaca-7B	
Methods		SGen ^{Semi} _{NoMS} (x)	SGen ^{Semi} (o)	SGen ^{Semi} _{NoMS} (x)	SGen ^{Semi} (o)
f_{M_1}	FDR-E	0.0609	0.1589	0.0359	0.0685
	efficiency	0.2829	0.7334	0.1580	0.3173
f_{M_2}	FDR-E	0.1785	0.1589	0.0698	0.0685
	efficiency	0.7835	0.7334	0.3200	0.3173
average efficiency		0.5347	0.7334	0.2390	0.3173

More in Our Paper

- How the **mislabeling error in pseudo-labeling** is controlled via **conformal prediction** & affects the FDR-E bound computation.
- Supervised-learning algorithm** as a special case of the proposed semi-supervised learning algorithm.
- Defining and choosing a “good” **scoring function** in terms of designing a single-threshold selection function.