

Selective Generation for Controllable Language Models

Minjae Lee^{1*}, Kyungmin Kim^{1*}, Taesoo Kim², Sangdon Park¹

¹POSTECH, ²Georgia Tech

*Equal contribution

★ NeurIPS 2024 Spotlight Paper



Selective Classification

Motivation

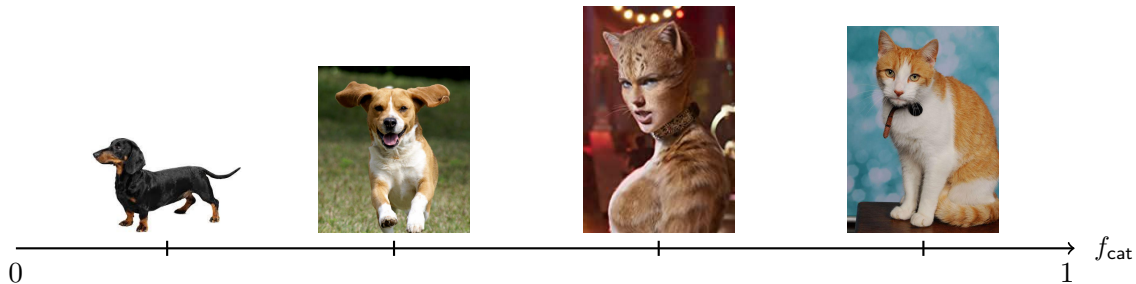
$$f_{\text{cat}} \left\{ \text{dog} \right\} < \tau \implies \text{Reject!}$$

- ▶ Selective classification is a **certified** risk control method, which rejects instances as needed, to grant a **desired risk** ε with **high probability** $1 - \delta$.

Why Selective Classification?

Motivation

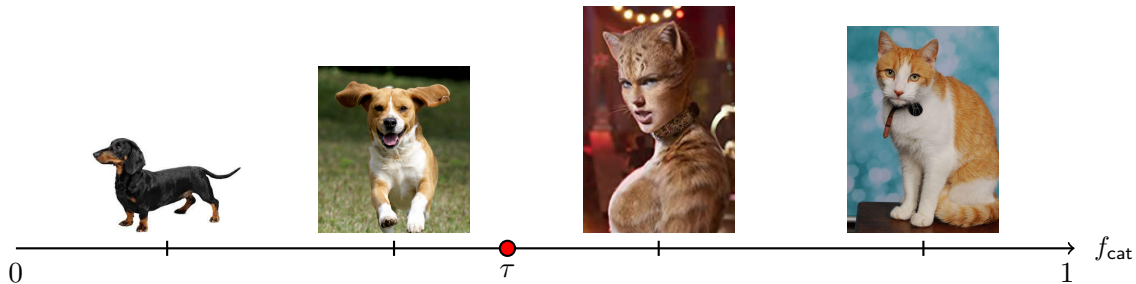
Training Phase



Why Selective Classification?

Motivation

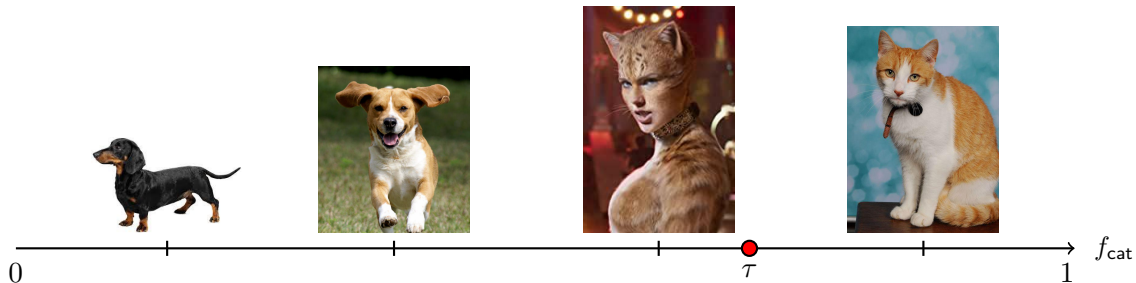
Training Phase



Why Selective Classification?

Motivation

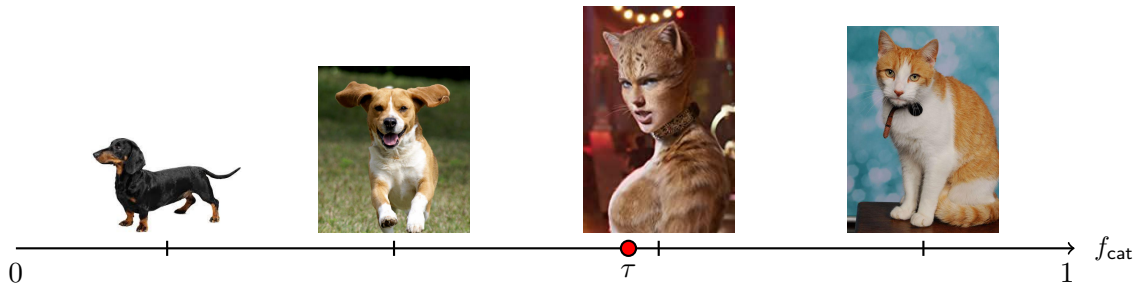
Training Phase



Why Selective Classification?

Motivation

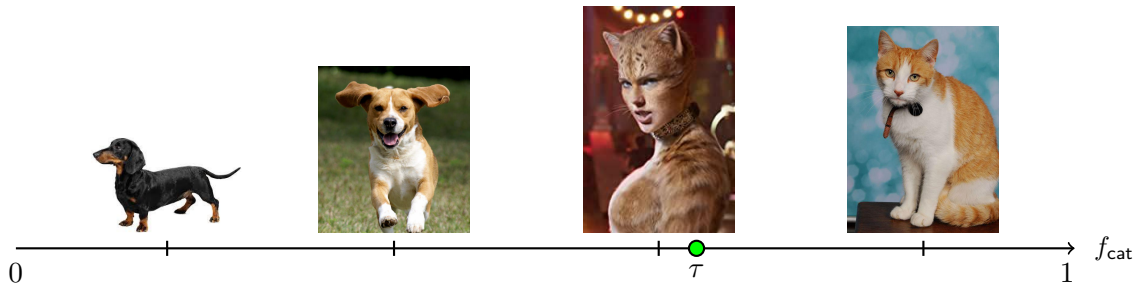
Training Phase



Why Selective Classification?

Motivation

Training Phase



Why Selective Classification?

Motivation

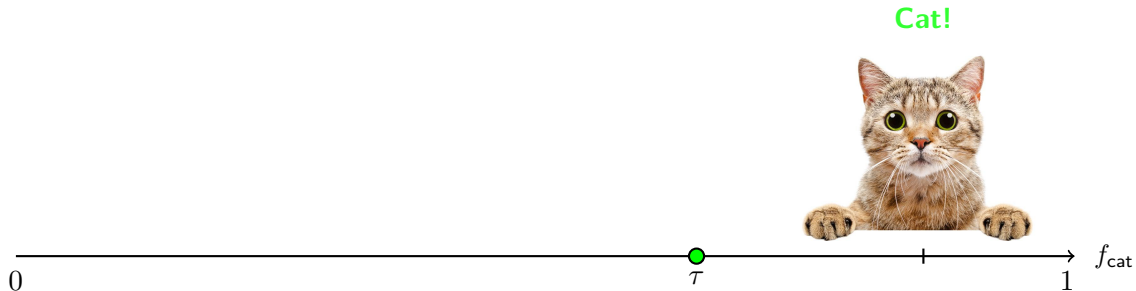
Inference Phase



Why Selective Classification?

Motivation

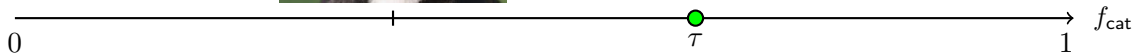
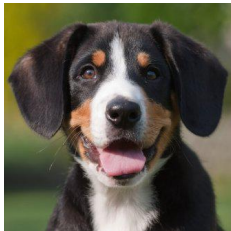
Inference Phase



Why Selective Classification?

Motivation

Inference Phase

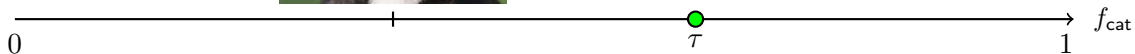
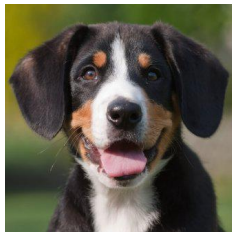


Why Selective Classification?

Motivation

Inference Phase

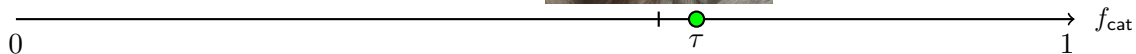
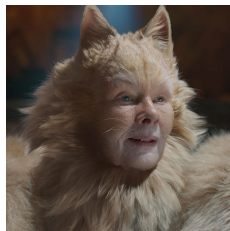
IDK...



Why Selective Classification?

Motivation

Inference Phase

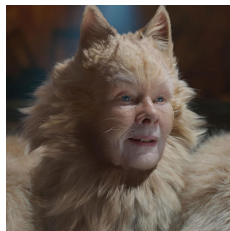


Why Selective Classification?

Motivation

Inference Phase

IDK...



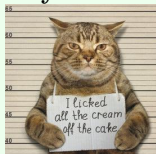
Why Textual Entailment?

Motivation

Classification

T : Classify the image below.

$\hat{y} = \text{cat}$



, $y = \text{cat}$

- ▶ **Selective prediction** is also important to be applied to **generative tasks**.
- ▶ However, unlike exact match (EM) in **classification**, it is difficult to define a **correctness metric**.

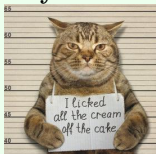
Why Textual Entailment?

Motivation

Classification

T : Classify the image below.

$\hat{y} = \text{cat}$



,

$y = \text{cat}$

$\hat{y} \stackrel{\text{EM!}}{=} y$

- ▶ **Selective prediction** is also important to be applied to **generative tasks**.
- ▶ However, unlike exact match (EM) in **classification**, it is difficult to define a **correctness metric**.

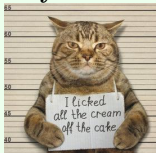
Why Textual Entailment?

Motivation

Classification

T : Classify the image below.

$\hat{y} = \text{cat}$



,

$y = \text{cat}$

$\hat{y} \stackrel{\text{EM!}}{=} y$

Generation

T : What is the objective of tour de france?

$\mathcal{Y} = \left\{ \begin{array}{l} \text{bike race} \\ \text{cycle race} \\ \text{biking competition} \\ \text{to pick the best cyclist} \\ \vdots \end{array} \right\}, y = \text{bicycle race}$

- ▶ **Selective prediction** is also important to be applied to **generative tasks**.
- ▶ However, unlike exact match (EM) in **classification**, it is difficult to define a **correctness metric**.

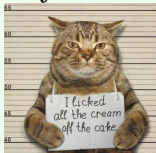
Why Textual Entailment?

Motivation

Classification

T : Classify the image below.

$\hat{y} = \text{cat}$



,

$y = \text{cat}$

$\hat{y} \stackrel{\text{EM!}}{=} y$

Generation

T : What is the objective of tour de france?

$\mathcal{Y} = \left\{ \begin{array}{l} \text{bike race} \\ \text{cycle race} \\ \text{biking competition} \\ \text{to pick the best cyclist} \\ \vdots \end{array} \right\}, y = \text{bicycle race}$

$\hat{y} \stackrel{\text{EM?}}{=} y$

- ▶ **Selective prediction** is also important to be applied to **generative tasks**.
- ▶ However, unlike exact match (EM) in **classification**, it is difficult to define a **correctness metric**.

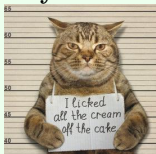
Why Textual Entailment?

Motivation

Classification

T : Classify the image below.

$\hat{y} = \text{cat}$



$y = \text{cat}$

$\hat{y} \stackrel{\text{EM!}}{=} y$

Generation

T : What is the objective of tour de france?

$\mathcal{Y} = \left\{ \begin{array}{l} \text{bike race} \\ \text{cycle race} \\ \text{biking competition} \\ \text{to pick the best cyclist} \\ \vdots \end{array} \right\}, y = \text{bicycle race}$

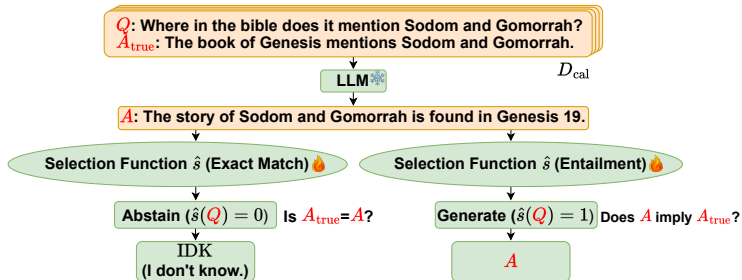
$\hat{y} \stackrel{\text{EM?}}{=} y$

- ▶ **Selective prediction** is also important to be applied to **generative tasks**.
- ▶ However, unlike exact match (EM) in **classification**, it is difficult to define a **correctness metric**.

\implies We employ **textual entailment**: $E_{\text{true}}(y) := \{\hat{y} \in \mathcal{Y} \mid \hat{y} \text{ implies } y\}$.

Why Semi-Supervised Learning?

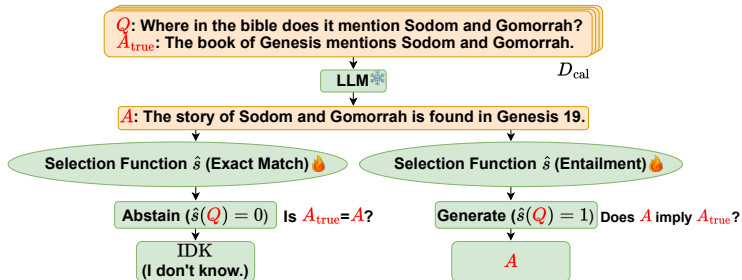
Motivation



- ▶ We can avoid *metric misalignment* in generation by leveraging entailment.
- ▶ However, labeling is expensive.

Why Semi-Supervised Learning?

Motivation

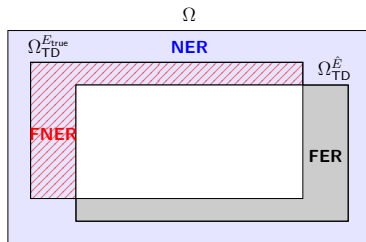


- ▶ We can avoid *metric misalignment* in generation by leveraging entailment.
- ▶ However, labeling is expensive.

⇒ We leverage question-answering pairs without entailment via **semi-supervised learning (SSL)**.

FDR-E Bound (1)

Method

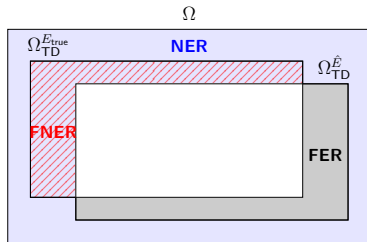


With the previously defined *textual entailment* $E_{\text{true}}(\mathbf{y})$, We can define **FDR-E**, the false discovery rate with respect to the textual entailment relation, as follows:

$$\mathbb{P}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y}) \mid \hat{S}(\mathbf{x}) \neq \text{IDK}\}$$

FDR-E Bound (2)

Method

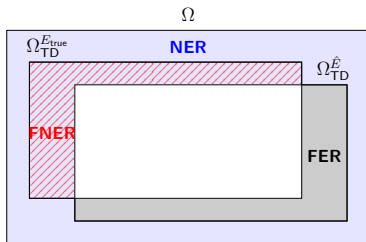


In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{(A)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 1\}}_{(B)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(C)} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 0\}}_{(D)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)}$$

FDR-E Bound (2)

Method

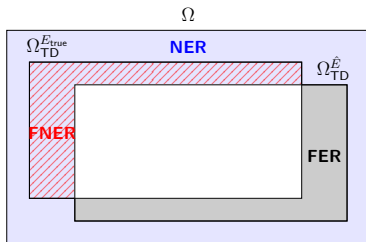


In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{(A)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 1\}}_{(B)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(C)} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 0\}}_{(D)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)}$$

FDR-E Bound (2)

Method

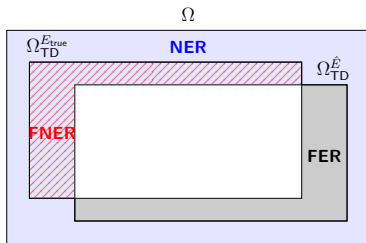


In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{(A)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 1\}}_{(B)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(C)} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 0\}}_{(D)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)}$$

FDR-E Bound (2)

Method



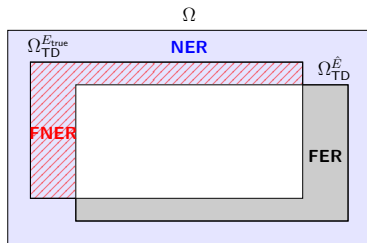
In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{(A)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 1\}}_{(B)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(C)} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 0\}}_{(D)}$$

$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)}$

FDR-E Bound (2)

Method



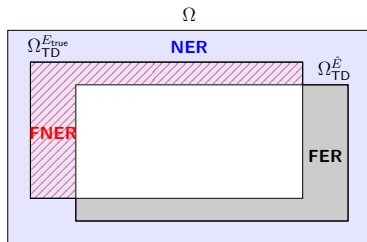
In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{(A)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 1\}}_{(B)} \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(C)} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{v = 0\}}_{(D)}$$

$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)}$

FDR-E Bound (2)

Method

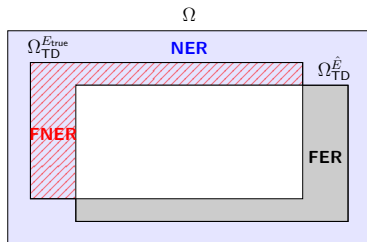


In the SSL setup, the FDR-E can be decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_S} \{C(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y})\}}_{\substack{\mathbb{P}_{\mathcal{D}_S} \{e = 0\} \\ \text{(E)}}} = \underbrace{\mathbb{P}_{\mathcal{D}_S} \{v = 1\}}_{\text{(B)}} \underbrace{\mathbb{P}_{\mathcal{D}_S} \{e = 0\}}_{\text{(C)}} + \underbrace{\mathbb{P}_{\mathcal{D}_S} \{v = 0\}}_{\text{(D)}}$$

FDR-E Bound (2)

Method



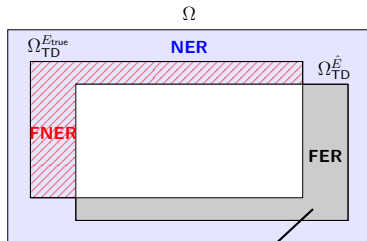
Lemma 1

(E) is decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0, \hat{e} = 1\}}_{FER} - \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 1, \hat{e} = 0\}}_{FNER} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{\hat{e} = 0\}}_{NER}.$$

FDR-E Bound (2)

Method



Lemma 1

(E) is decomposed as follows:

$$\underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0\}}_{(E)} = \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 0, \hat{e} = 1\}}_{\text{FER}} - \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{e = 1, \hat{e} = 0\}}_{\text{FNER}} + \underbrace{\mathbb{P}_{\mathcal{D}_{\hat{S}}}\{\hat{e} = 0\}}_{\text{NER}}.$$

$\leq \epsilon_E$

FDR-E Bound (3)

Method

Entailment Set Learning

\mathcal{A}_{FER} returns \hat{E} which controls the **FER** of pseudo-labeled examples, *i.e.*,

$$\mathbb{P}\{\mathcal{R}_{\text{FER}}(\hat{E}) \leq \varepsilon_E\} \geq 1 - \delta_E.$$

- ▶ We use \hat{E} as a pseudo-labeling function for SSL – see our paper!

Lemma 2

If $\hat{E} := \mathcal{A}_{\text{FER}}(\hat{\mathbf{Z}}_E)$ satisfies the above guarantee, we have

$$\mathbb{P}_{\mathcal{D}}\{e = 0\} \leq \varepsilon_E - L_{\text{Binom}}(\hat{k}; |\hat{\mathbf{Z}}_E|, \delta'_E/2) + U_{\text{Binom}}(\hat{l}; |\hat{\mathbf{Z}}_U|, \delta'_S) =: U_{\text{SSL}}$$

- ▶ We find an optimal ε_E that minimizes U_{SSL} , resulting $U_{\text{SSL}}^{\text{OPT}}$ – see our paper!

Controllable Guarantee

Method

Algorithm

Our **semi-supervised** method $\mathcal{A}_{\text{SGen}}^{\text{Semi}}$ solves the following optimization problem:

$$\text{find}_{\hat{S} \in \mathcal{H}} \hat{S} \quad \text{subj. to} \quad w_{\text{SL}} U_{\text{SL}} + w_{\text{SSL}} U_{\text{SSL}}^{\text{OPT}} \leq \varepsilon_S$$

Theorem 1

$\mathcal{A}_{\text{SGen}}^{\text{Semi}}$ satisfies the following controllable guarantee on the FDR-E, i.e.,

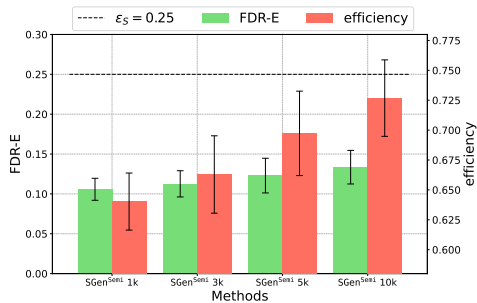
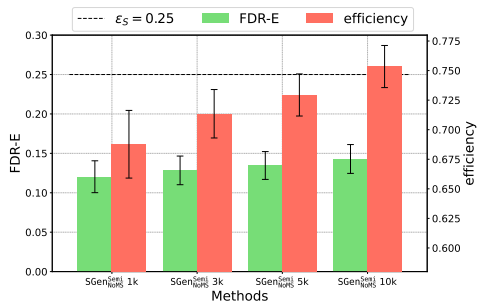
$$\mathbb{P} \left\{ \mathbb{P} \{ G(\mathbf{x}) \notin E_{\text{true}}(\mathbf{y}) \mid \hat{S}(\mathbf{x}) \neq \text{IDK} \} \leq \hat{U} \right\} \geq 1 - \delta.$$

Experiment

Question x	Who is the actor who plays Draco Malfoy?	When did the movie Benjamin Button come out?
Correct Answer y	Thomas Andrew Felton plays Draco Malfoy in the Harry Potter movies.	The movie Benjamin Button come out December 25, 2008
Generated Answer $G(x)$	The actor who plays Draco Malfoy is Tom Felton. (correct)	The movie The Curious Journey of Benjamin Button was released in 2008. (correct)
$S\text{Gen}_{EM}$	rejected	rejected
$S\text{Gen}^{\text{Semi}}$ (ours)	accepted	accepted

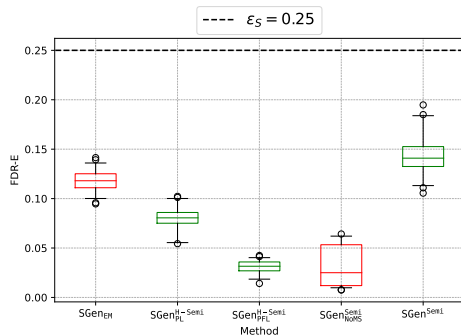
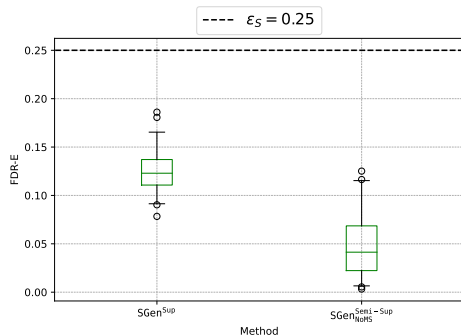
- ▶ $S\text{Gen}^{\text{Semi}}$ can capture correctness better than $S\text{Gen}_{EM}$.

Experiment



- More unlabeled samples are beneficial to achieving better efficiency.

Experiment



- The FDR-E for \hat{S} is well controlled below ϵ_S , desired FDR-E, under the test environment.

Conclusion

- ▶ We leverage logical **entailment** and propose a novel **semi-supervised** learning approach for **selective generation**, demonstrating its *theoretical* and *empirical* efficacy.

