# Commonsense Reasoning with PTLMs

*What do you fill with ink to write notes on a piece of copy paper ?*

*(A) Fountain pen*
*(B) Pencil case*
*(C) Printer*
*(D) Notepad*

# Commonsense Reasoning with PTLMs

*What do you fill with ink to write notes on a piece of copy paper ?*

**(A) Fountain pen**
(B)  Pencil case
(C)  Printer
(D)  Notepad

# Commonsense Reasoning with PTLMs

*What do you fill with ink to write notes on a piece of copy paper ?*

**(A) Fountain pen**
(B) Pencil case
(C) Printer
(D) Notepad



**AI2** Allen Institute for AI

🌐 **UNIFIED-QA**

Enter a question to see what answer our UnifiedQA gives. You can also use one of the examples below.

Examples:

Format: Multiple-Choice (Elementary-school-science)-1

Raw string input (first comes the question, then your list your candidates or paragraph; use "\n" as separator):

What do you fill with ink to write notes on a piece of copy paper? \n(A) fountain pen (B) pencil case (C) printer (D) notepad

Submit

**Input:**

What do you fill with ink to write notes on a piece of copy paper?

(A) fountain pen (B) pencil case (C) printer (D) notepad

**Output:**

Prediction [small, 60 million parameters]: pencil case

Prediction [large, 770 million parameters]: printer

*Base : pencil case*
*Large : printer*

Fails to reason with the **concept-centric knowledge**

# Current PTLMs

**PTLMs**

**Corpus**

… **Copy paper** is thinner than **printer** paper, which doesn't make a huge difference when you're printing text, but it does when you're printing large images. Images require a lot of **ink** and because **copy paper** has a thinner structure, the **ink** will need to spread out more for the **paper** to absorb it all. ...

**Pre-train**

Text Infilling / MLM

**AI2** Allen Institute for AI

🌐 **UNIFIED-QA**

Enter a question to see what answer our UnifiedQA gives. You can also use one of the examples below.

Examples:

Format: Multiple-Choice (Elementary-school-science)-1

Raw string input (first comes the question, then your list your candidates or paragraph; use "\n" as separator):

What do you fill with ink to write notes on a piece of copy paper? \n(A) fountain pen (B) pencil case (C) printer (D) notepad

Submit

**Input:**

What do you fill with ink to write notes on a piece of copy paper?

(A) fountain pen (B) pencil case (C) printer (D) notepad

**Output:**

Prediction [small, 60 million parameters]: pencil case

Prediction [large, 770 million parameters]: printer

*Base : pencil case*
*Large : printer*

The model may be **sensitive** to the **co-occurence (ink, copy, paper)**

5

# How can we teach PTLMs to write and reason with Common sense Concepts ?

*What do you **fill** with **ink** to **write** notes on a piece of **copy paper** ?*

(A) **Fountain pen**
(B) Pencil case
(C) Printer
(D) Notepad

*fill, ink, write, copy paper*

*Fountain Pen*

# **Our idea** : Novel Self-supervised Objectives to improve common sense reasoning ability.

*Generate a sentence with the following concepts :*
*Hold Woman Position*

**Text-to-Text Transformer**

She was the first woman to hold the position.

**Generative Objective** : **Concept-to-Sentence Generation (C2S)**
Ask model to recover the original sentence
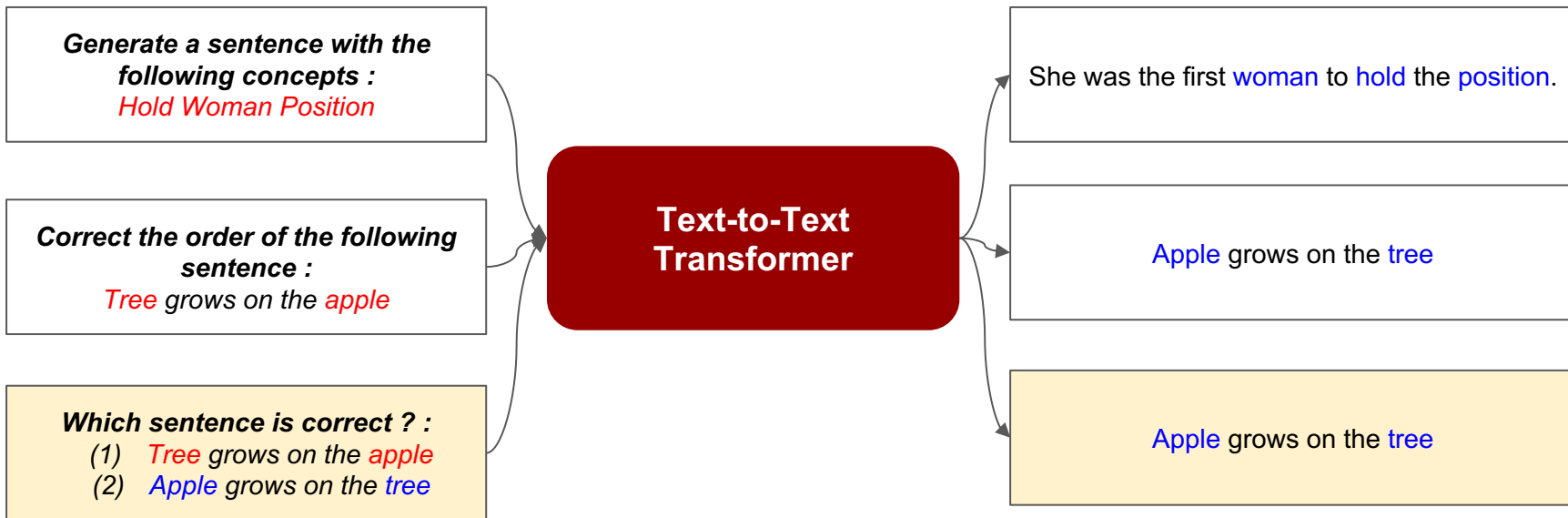given only a few unordered keywords of the sentence.

# **Our idea** : Novel Self-supervised Objectives to improve common sense reasoning ability.



| | | |
|---|---|---|
| ***Generate a sentence with the following concepts :***<br>*Hold Woman Position* | **Text-to-Text Transformer** | She was the first woman to hold the position. |
| ***Correct the order of the following sentence :***<br>*Tree grows on the apple* | | Apple grows on the tree |

**Generative Objective** : **Concept Order Recovering (COR)**
Ask model to recover the original sentence
given order-of-concept shuffled sentence.

# **Our idea** : Novel Self-supervised Objectives to improve common sense reasoning ability.



**Discriminative Objective** : **Generative QA**
Ask model to distinguish the real sentence from a concept-distracted sentence.

# CALM : Concept-Aware Language Model

Original Sentence *x*

She was the first woman to hold the position.

**Extract Concept Set** *C*
(*woman, hold, position*)

(1) Given an input sentence *x* (*"She was the first woman to hold the position."*),
extract concept-set *C* (*woman, hold, position*).

# CALM : **C**oncept-**A**ware **L**anguage **M**odel

Original Sentence *x*

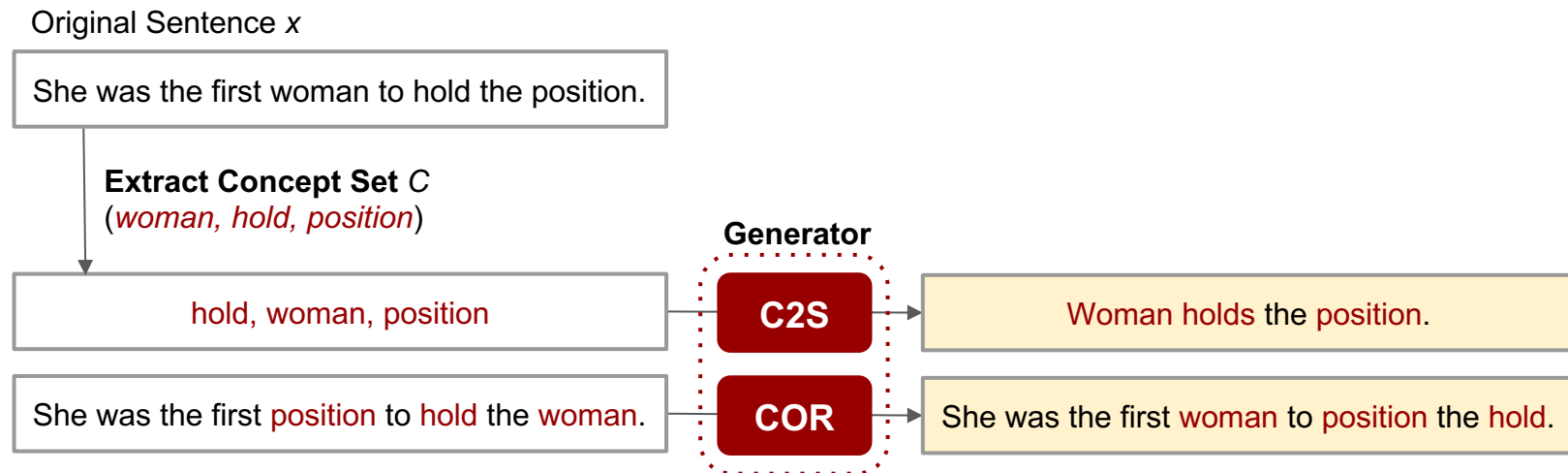> She was the first woman to hold the position.

**Extract Concept Set** *C*
(*woman, hold, position*)

> hold, woman, position

> She was the first position to hold the woman.

(1)  Given an input sentence *x* (*"She was the first woman to hold the position."*),
     extract concept-set *C* (*woman, hold, position*).

(1)  Given *x* and *C*, produce corrupted source sentence *x'* either for **C2S** and **COR**

# CALM : **C**oncept-**A**ware **L**anguage **M**odel

Original Sentence *x*

She was the first woman to hold the position.

**Extract Concept Set** *C*
(*woman, hold, position*)

**Generator**

| hold, woman, position | **C2S** | Woman holds the position. |

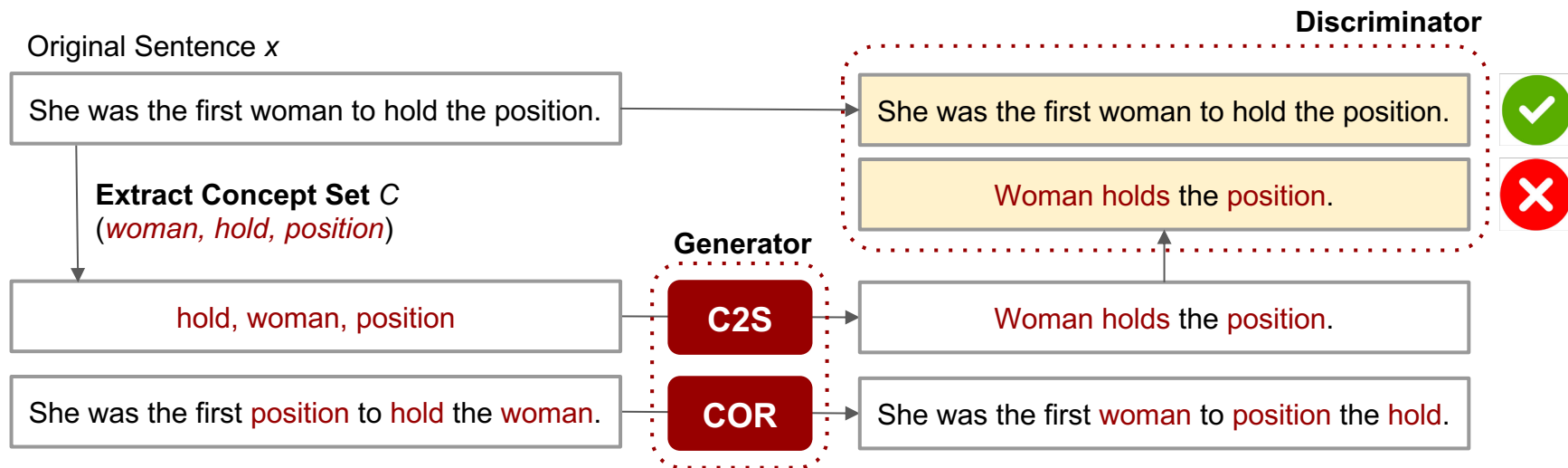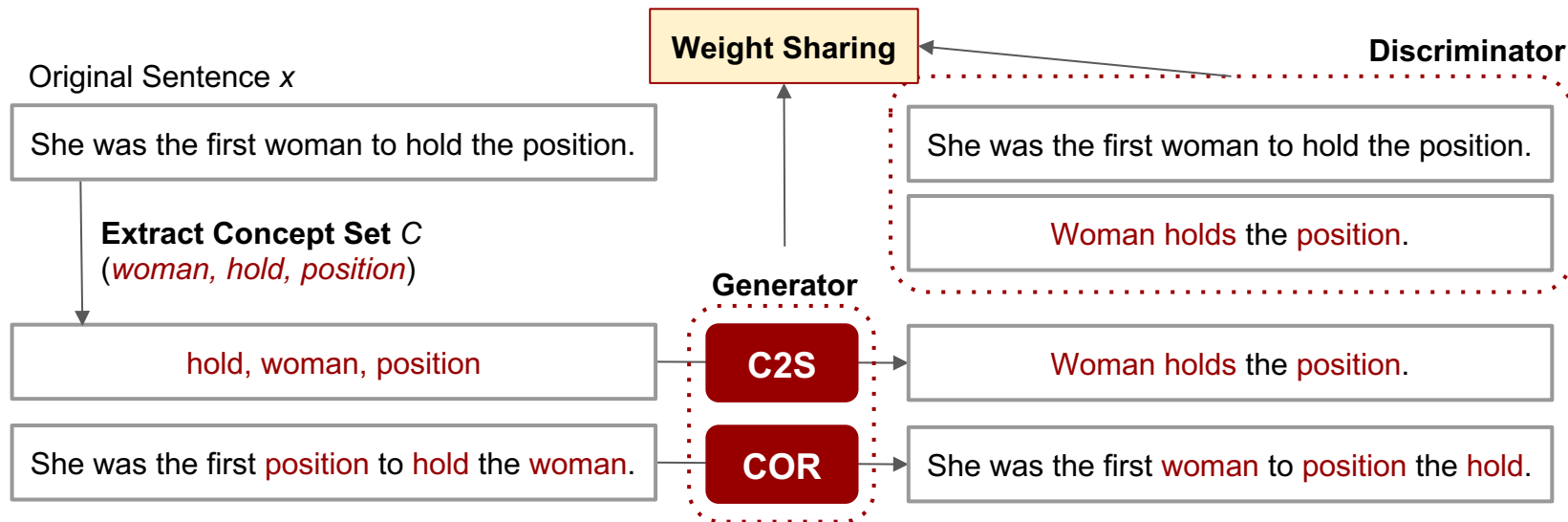| She was the first position to hold the woman. | **COR** | She was the first woman to position the hold. |

(1) Given an input sentence *x* (*"She was the first woman to hold the position."*),
    extract concept-set *C* (*woman, hold, position*).
(1) Given *x* and *C*, produce corrupted source sentence *x'* either for **C2S** and **COR**
(2) The **generator** trained with **C2S** and **COR** recovers sentence *x'* to distractor *x"*

# CALM : **C**oncept-**A**ware **L**anguage **M**odel

Original Sentence *x*

She was the first woman to hold the position.

**Discriminator**

She was the first woman to hold the position. ✅

Woman holds the position. ❌

**Extract Concept Set** *C*
(*woman, hold, position*)

**Generator**

hold, woman, position

**C2S** → Woman holds the position.

She was the first position to hold the woman.

**COR** → She was the first woman to position the hold.

(1) Given an input sentence *x* ("*She was the first woman to hold the position.*"),
    extract concept-set *C* (*woman, hold, position*).
(1) Given *x* and *C*, produce corrupted source sentence *x'* either for **C2S** and **COR**
(2) The **generator** trained with **C2S** and **COR** recovers sentence *x'* to distractor *x"*
(3) The **discriminator** is trained to distinguish truth sentence from distractor *x"*

13

# CALM : **C**oncept-**A**ware **L**anguage **M**odel

**Weight Sharing**

**Discriminator**

Original Sentence *x*

| She was the first woman to hold the position. |

**Extract Concept Set** *C*
(*woman, hold, position*)

**Generator**

| She was the first woman to hold the position. |

| Woman holds the position. |

| hold, woman, position | → | **C2S** | → | Woman holds the position. |

| She was the first position to hold the woman. | | **COR** | → | She was the first woman to position the hold. |

---

(1) Given an input sentence *x* ("*She was the first woman to hold the position.*"),
    extract concept-set *C* (*woman, hold, position*).
(1) Given *x* and *C*, produce corrupted source sentence *x'* either for **C2S** and **COR**
(2) The **generator** trained with **C2S** and **COR** recovers sentence *x'* to distractor *x"*
(3) The **discriminator** is trained to distinguish truth sentence from distractor *x"*

# Is **CALM** **reason** with **concepts** ? *Yes !*

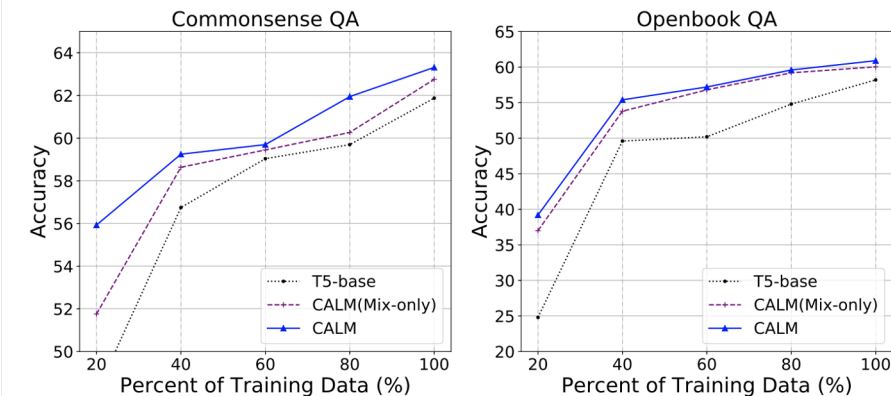| Methods | CSQA | OBQA | PIQA | aNLI |
|---|---|---|---|---|
| | Accuracy | | | |
| T5-base | 61.88($\pm$0.08) | 58.20($\pm$1.0) | 68.14($\pm$0.73) | 61.10($\pm$0.38) |
| T5-base w/ additional epochs | 61.92($\pm$0.45) | 58.10($\pm$0.9) | 68.19($\pm$0.77) | 61.15($\pm$0.52) |
| T5-base + SSM | 62.08($\pm$0.41) | 58.30($\pm$0.8) | 68.27($\pm$0.71) | 61.25($\pm$0.51) |
| CALM (Generative-Only) | 62.28($\pm$0.36) | 58.90($\pm$0.4) | 68.91($\pm$0.88) | 60.95($\pm$0.46) |
| CALM (Contrastive-Only) | 62.73($\pm$0.41) | 59.30($\pm$0.3) | <u>70.67</u>($\pm$0.98) | 61.35($\pm$0.06) |
| CALM (Mix-only) | <u>63.02</u>($\pm$0.47) | <u>60.40</u>($\pm$0.4) | 70.07($\pm$0.98) | <u>62.79</u>($\pm$0.55) |
| CALM (w/o Mix warmup) | 62.18($\pm$0.48) | 59.00($\pm$0.5) | 69.21($\pm$0.57) | 61.25($\pm$0.55) |
| CALM | **63.32($\pm$0.35)** | **60.90($\pm$0.4)** | **71.01($\pm$0.61)** | **63.20($\pm$0.52)** |

Experimental results on commonsense reasoning dataset.

**CALM** consistently and significantly **<u>outperforms</u>**
the backbone T5-base model.

# Is **CALM** **reason** with **concepts** ? *Yes !*

| Methods | CSQA | OBQA | PIQA | aNLI |
|---|---|---|---|---|
| | Accuracy (official dev) | | | |
| BERT-large | 57.06($\pm$0.12) | 60.40($\pm$0.6) | 67.08($\pm$0.61) | 66.75($\pm$0.61) |
| T5-large | 69.81($\pm$1.02) | 61.40($\pm$1.0) | 72.19($\pm$1.09) | 75.54($\pm$1.22) |
| CALM-large (Mix-only) | 70.26($\pm$0.23) | 62.50($\pm$1.0) | 73.70($\pm$1.09) | 75.99($\pm$1.26) |
| CALM-large | 71.31($\pm$0.04) | **66.00($\pm$1.0)** | 75.11($\pm$1.65) | 77.12($\pm$0.34) |

Effective in Large Models.



Performance of compared models
fine-tuned with <u>different fraction</u> of the dataset

## Performance is **consistent** in
## <u>large</u> model & <u>**different fraction**</u> of the dataset.

16

# Is **CALM** **write** with **concepts** ? *Yes !*

| Methods | Params | CommonGEN | | | |
|---|---|---|---|---|---|
| | | BLEU-4 | METEOR | CIDEr | SPICE |
| GPT-2 (Radford et al., 2019) | 774M | 21.10 | 26.20 | 12.15 | 25.90 |
| UniLM (Dong et al., 2019) | 340M | 27.70 | 29.70 | 14.85 | 30.20 |
| BART (Lewis et al., 2020) | 406M | 26.30 | 30.90 | 13.92 | 30.60 |
| T5-Base (Raffel et al., 2019) | 220M | 16.40 | 23.00 | 9.16 | 22.00 |
| T5-Large (Raffel et al., 2019) | 770M | 28.60 | 30.10 | 14.96 | 31.60 |
| KG-BART (Liu et al., 2020) | 406M | **30.90** | **32.40** | **16.83** | 32.70 |
| Our T5-Base | 220M | 24.90 | 31.20 | 12.99 | 32.40 |
| CALM | 220M | 26.40 | 31.40 | 13.88 | **33.00** |

(Left) : Comparison between PTLMs
(Below) : Comparison on generated sentences with same concept-set

| Concept-set | T5-base | CALM |
|---|---|---|
| Grass, Dog, Ball, Chase | a dog is chased by a ball on the grass. | dog chasing a ball in the grass. |
| Net, Cast, Boat, Water | fishing boat casts a net in the water. | fisherman casts a net into the water from a fishing boat. |
| Hole, Tree, Plant, Dig | a man digs a hole in a tree to plant a new tree . he digs the | man digging a hole to plant a tree. |
| Ingredient, Add, Pan, Fry | a pan filled with ingredients adds a touch of spice to the fry . | add the ingredients to a pan and fry. |
| Water, Hold, Hand, Walk | A man holding a hand and walking in the water. A man is holding water. | man holding a bottle of water in his hand as he walks down the street. |
| Place, Use, Metal tool | A man uses a metal tool to make a piece of metal. | woman uses a metal tool to make a piece of jewelry. |
| Hair, Wax, Apply, Remove | remove the wax from your hair and apply it to your hair . | woman applies wax to her hair and then removes it with a comb. |
| Sidewalk, Dog, Walk, Leash | A dog walking on a leash on the sidewalk. | dog walking on a sidewalk with a leash. |

# Summary

- **Novel self-supervised strategies** for concept-centric Common Sense
  - Concept to Sentence
  - Concept Order Recovering
  - Generative QA

- **Two-stage training strategy**
  - Generator and Discriminator

Text-to-Text models can be pre-trained with **better parameter** and **sample efficiency** by carefully designed **self-supervised objectives** that focus on the ability required by target tasks.