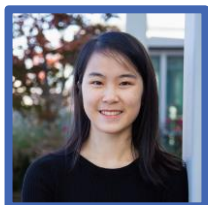


On the Influence of Masking Policies in Intermediate Pre-training



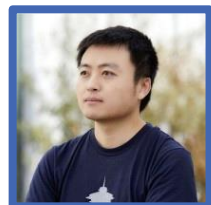
Qinyuan Ye^{1,†}

† Work partially done while at Facebook AI



Belinda Z. Li^{2,†}

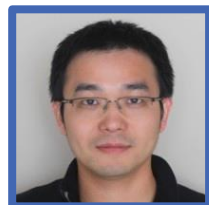
† Work partially done while at Facebook AI



Sinong Wang³



Benjamin Bolte³



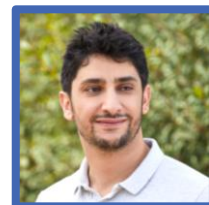
Hao Ma³



Wen-tau Yih³



Xiang Ren¹



Madian Khabsa³

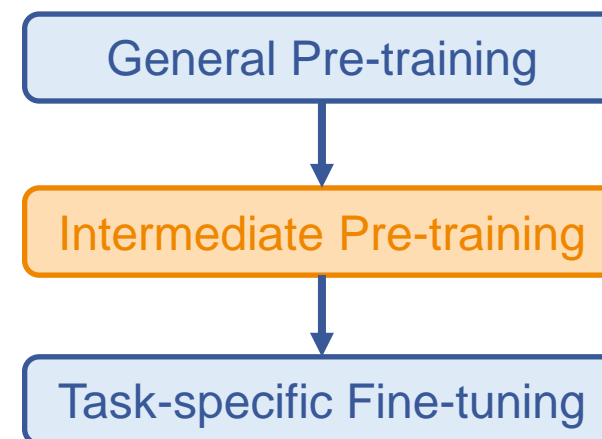
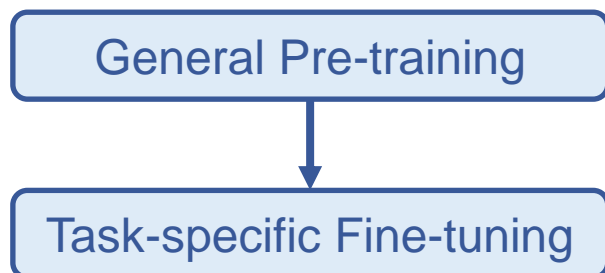
1 **USC**
Viterbi
*Department of
Computer Science*

2 
MIT CSAIL

3 **facebook** Artificial Intelligence

Motivation

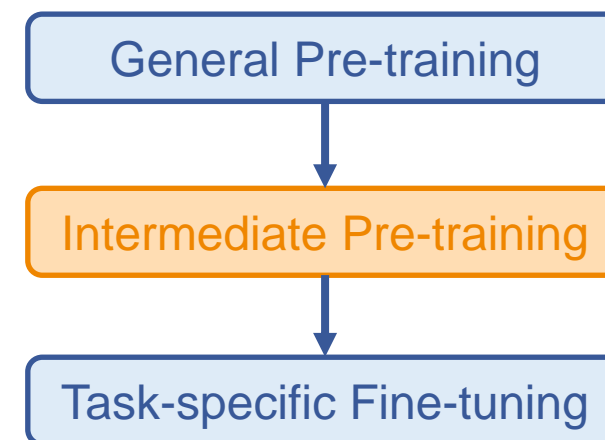
- “*Pre-train then fine-tune*” is a predominant pipeline in NLP.
- Inserting an *intermediate stage with a hand-crafted masking policy* can be helpful.



masking named entities → better closed-book QA models
(Roberts et al., 2020)

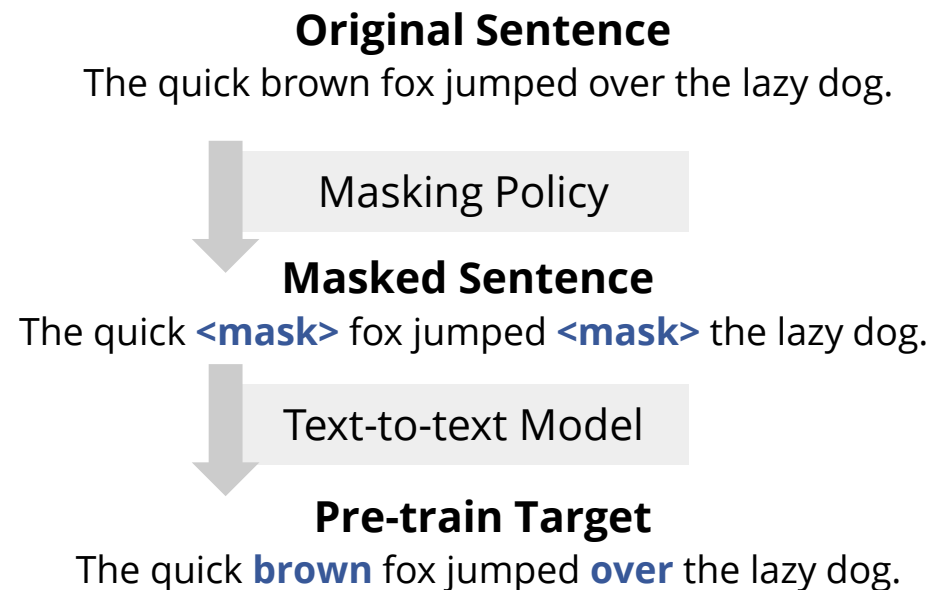
Motivation

- In this work we offer a large-scale **empirical study** to investigate **the influence of masking policies in intermediate pre-training**.
- We aim to answer
 - **In what cases** such intermediate pre-training is helpful
 - Whether hand-crafted heuristic objectives are **near-optimal**
 - Whether a masking policy designed for one task is **generalizable** beyond that task



Preliminaries – Masked Language Modeling

- Masked language modeling (MLM) and variants are common for pre-training large-scale transformers, e.g., BERT, RoBERTa, BART, T5.



Analysis Setup – Training Pipeline

General Pre-training with Random Masks

Source: In Newtonian physics, free fall <mask> motion of a <mask> where gravity is <mask> only <mask> upon it.

Text-to-text
Model

Target: In Newtonian physics, free fall is any motion of a body where gravity is the only force acting upon it.

Intermediate Pre-training

Source: In Newtonian physics, free fall is any motion of a body where <mask> is the only force acting upon it.

Text-to-text
Model

Target: gravity

Task-specific Fine-tuning

Source: which force acts on an object in free fall

Text-to-text
Model

Target: gravity

We ensure that the masking policy is the only variable in this pipeline.

Analysis Setup – Downstream Tasks

Closed-book QA

TQA: TriviaQA
WQ: WebQuestions
NQ: Natural Questions

Knowledge-intensive Tasks

AY2: Entity Linking
ZSRE: Zero-shot Relation Extraction
WoW: Dialogue

Multiple-choice QA

WIQA
QuaRTz
ROPES

Each example is a source-target pair (s, t) , accompanied with a context paragraph c

Compared Masking Policies (1/3)

Heuristic, Supervised, Meta-learned

- **Heuristics**

Original Objective (+Orig)

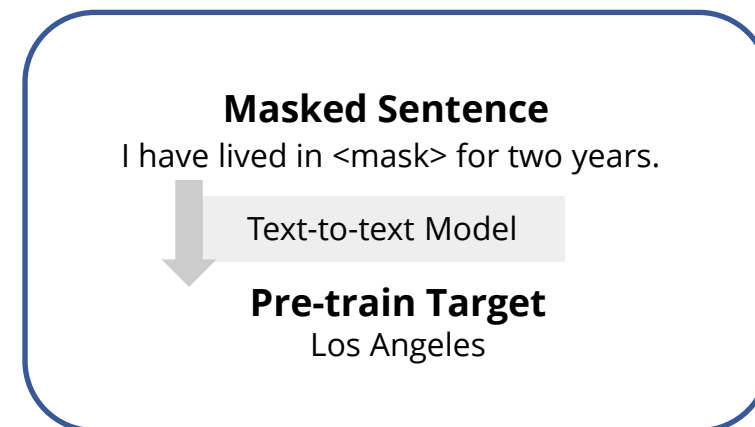
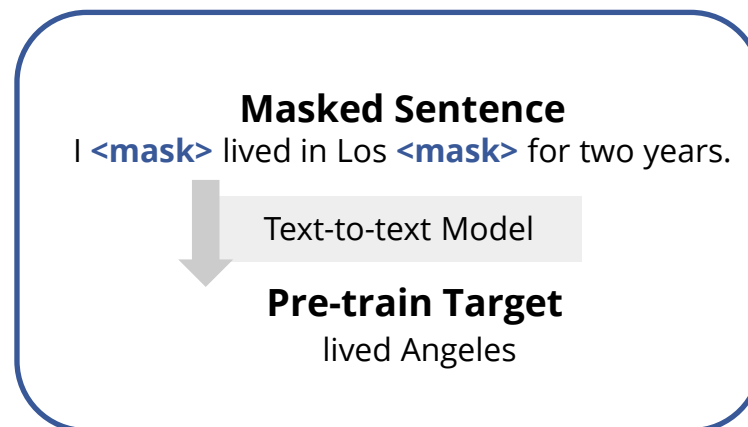
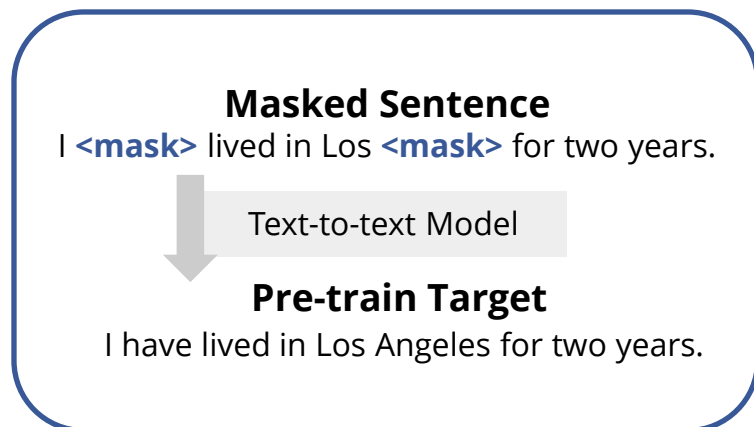
Mask 15% randomly selected tokens and recover full sequence

Random Masking (+Rand)

Mask 15% randomly selected tokens and recover the masked tokens

Salient Span Masking (+SSM)

Mask and recover one named entity



Compared Masking Policies (2/3)

Heuristic, Supervised, Meta-learned

- **Supervised**

- Identify *likely-answers* and mask them
- The masking policies is similar to an extractive reading comprehension model
- The policy is trained with (*context, answer*) examples, *without the question*

Example

Input: [Charles, Schulz, was, the, creator, of, Snoopy]

Masking Policy

Start Index: [1,0,0,0,0,0,0]

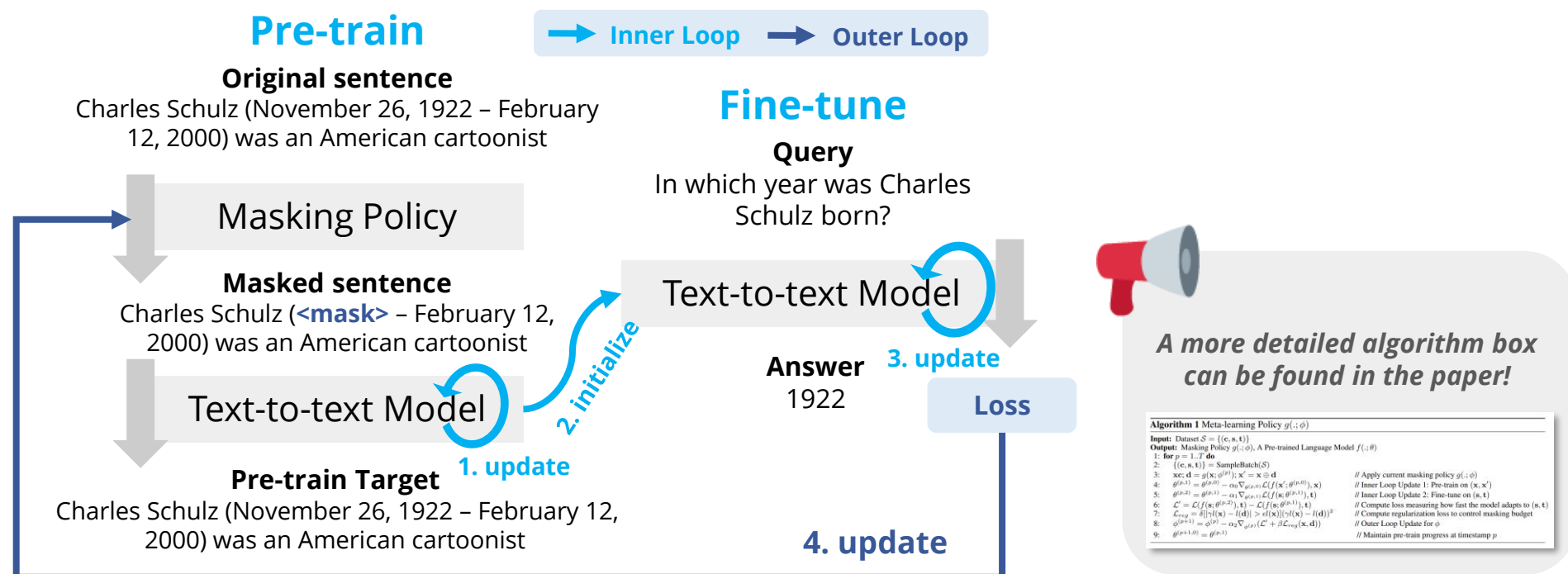
End Index: [0,1,0,0,0,0,0]

Compared Masking Policies (3/3)

Heuristic, Supervised, Meta-learned

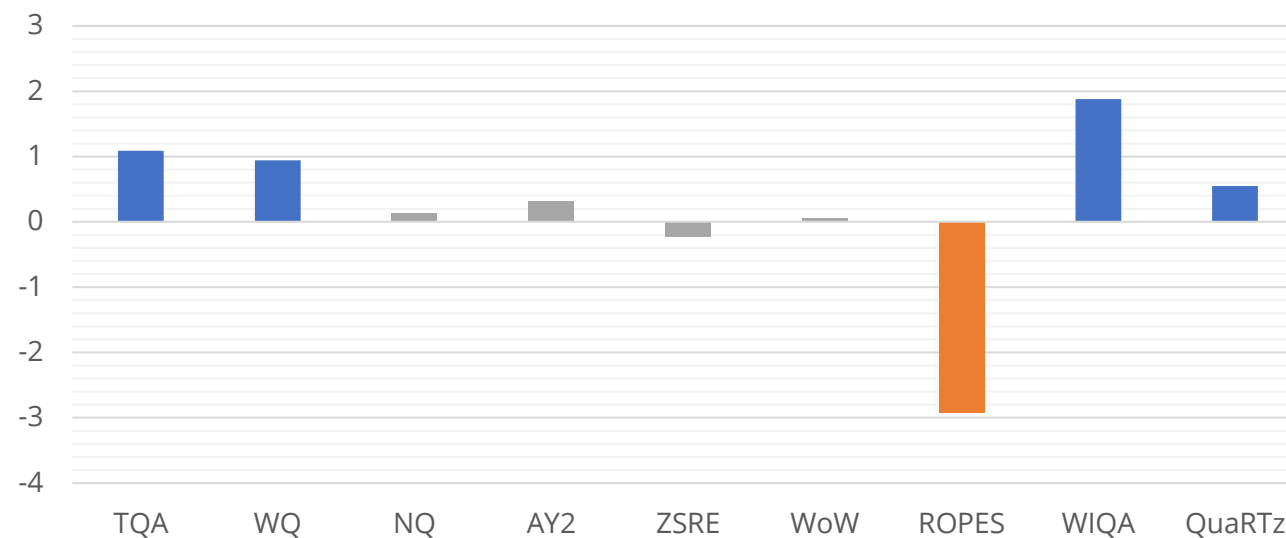
- **Meta-learned**

- Our goal is closely related to the concept of “*learning to learn*” (Schmidhuber, 1987; Thrun and Pratt, 1998).
- The masking policy should help the text-to-text to learn quickly when fine-tuned on downstream tasks.



Analysis – Comparing Heuristic Policies

Performance Improvement with +Orig

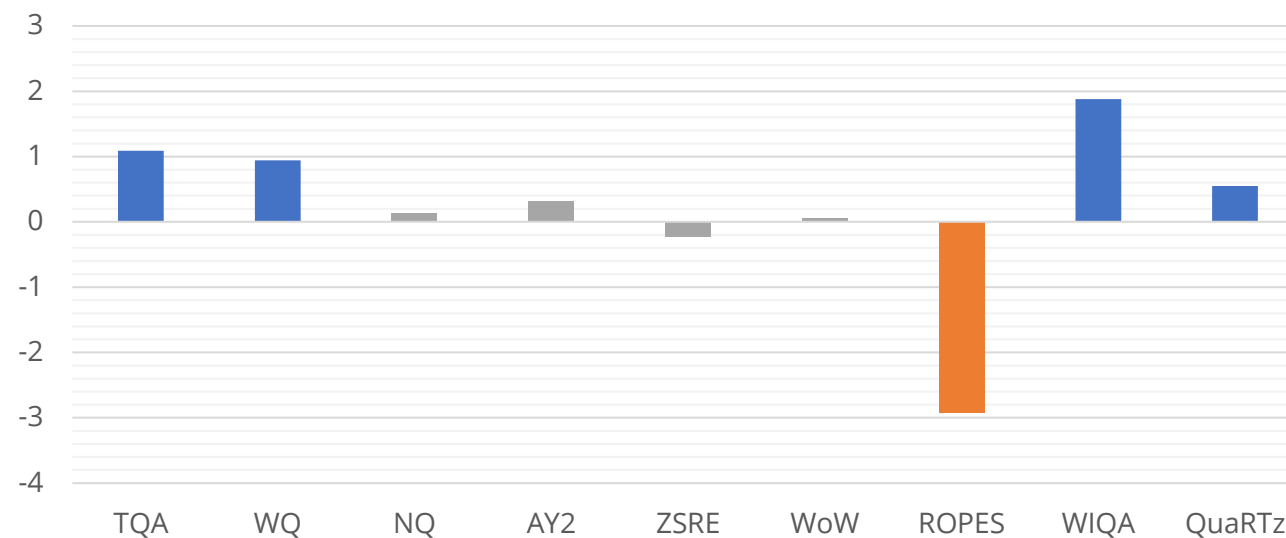


↑ x4 → x4 ↓ x1

Overall, intermediate pre-training with the original objective leads to improved performance.

Analysis – Comparing Heuristic Policies

Performance Improvement with +Orig

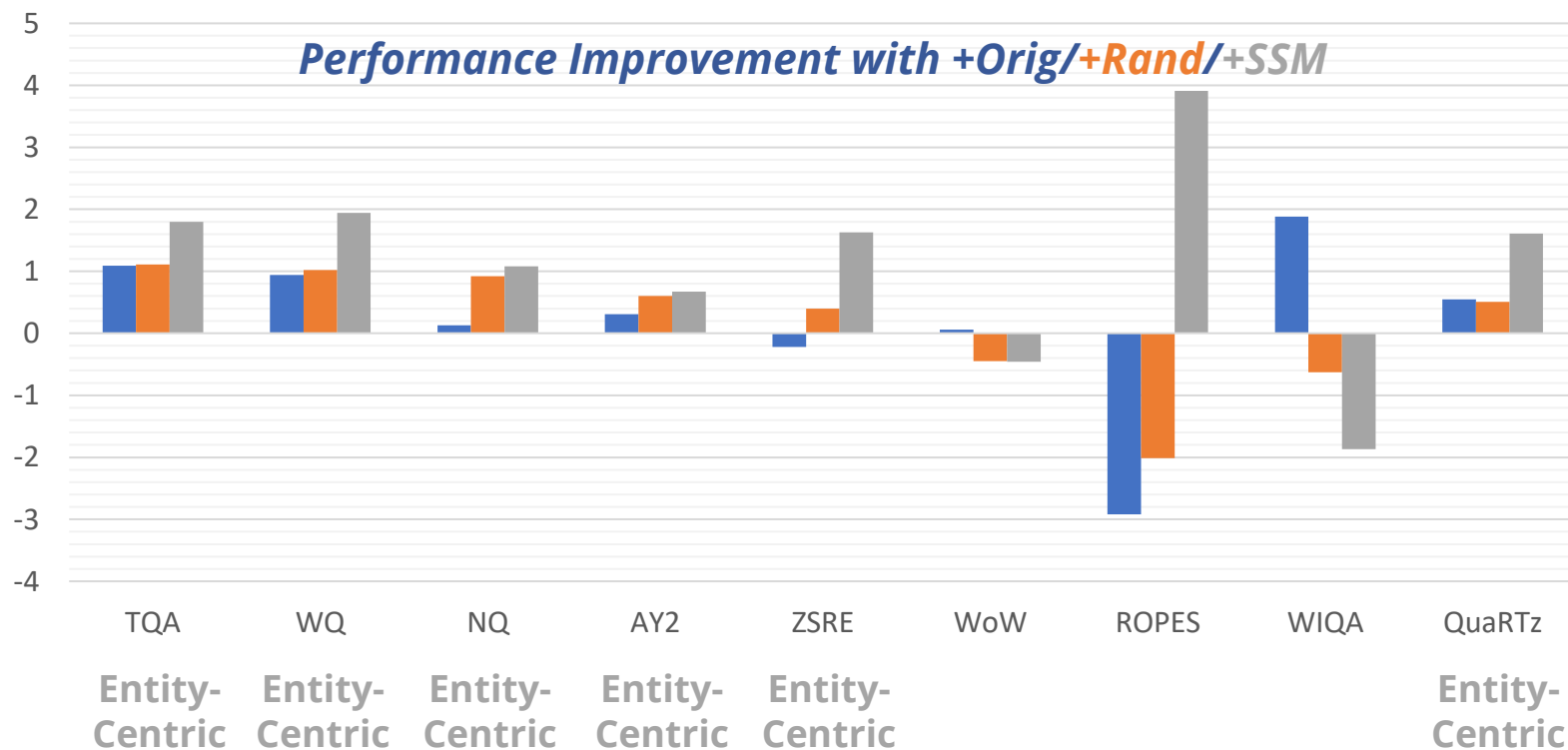


↑ x4 → x4 ↓ x1

ROPES dataset is *the only exception*. We found that intermediate pre-training may lead to *catastrophic forgetting*.

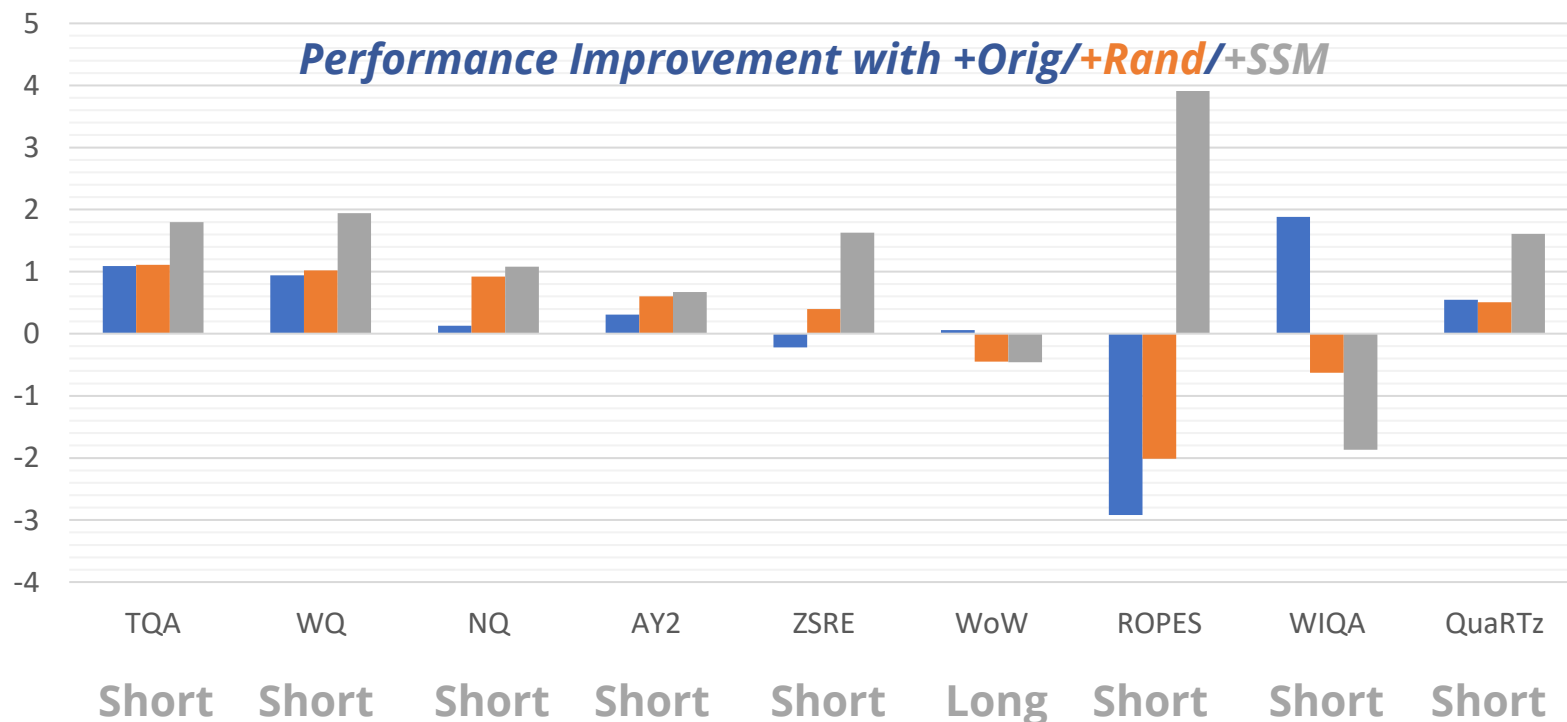
Pay attention to the corpus from which the dataset is created!

Analysis – Comparing Heuristic Policies

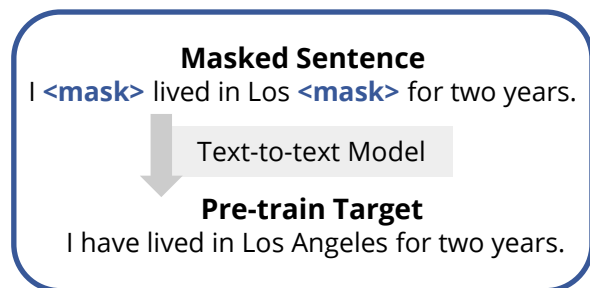


- SSM is beneficial for all entity-centric tasks.
- Use heuristic masking policies that resemble the downstream tasks, or masking information known to be important for the downstream task, tend to be helpful.

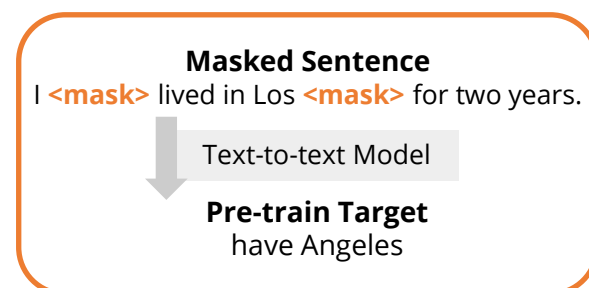
Analysis – Comparing Heuristic Policies



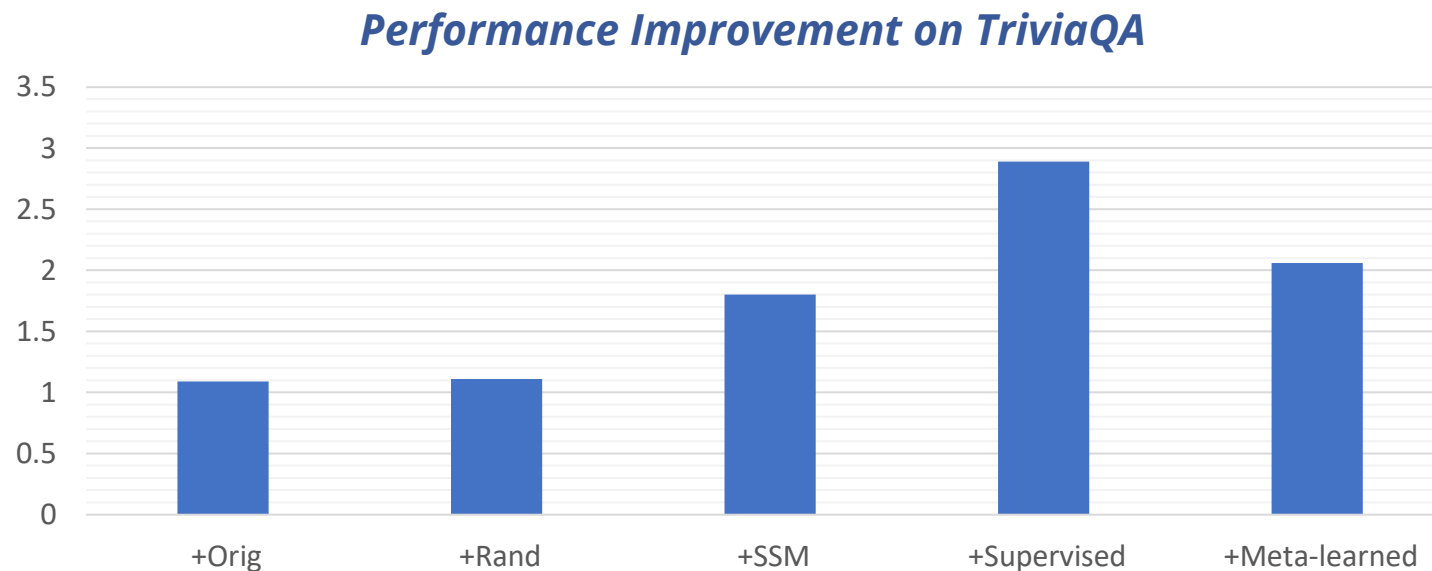
+Orig is better for long outputs



+Rand is better for short outputs



Analysis – Bringing in Learned Policies



Learned policies are most successful on TriviaQA.

Meta-learned policies also outperform +Orig on NQ (+0.98 EM), ZSRE (+0.32 EM) and ROPES (+10.03 Acc.)

Analysis – Bringing in Learned Policies

Additional Observations

Improved Learning Efficiency: Supervised policy has better learning efficiency than SSM on TriviaQA.

Generalization of Learned Policies: Learned masking policies *can positively transfer*. However more investigation is needed.

Overfitting on ZSRE: The learned policy may *overfit to the training data* if there is a mismatch between train/test data

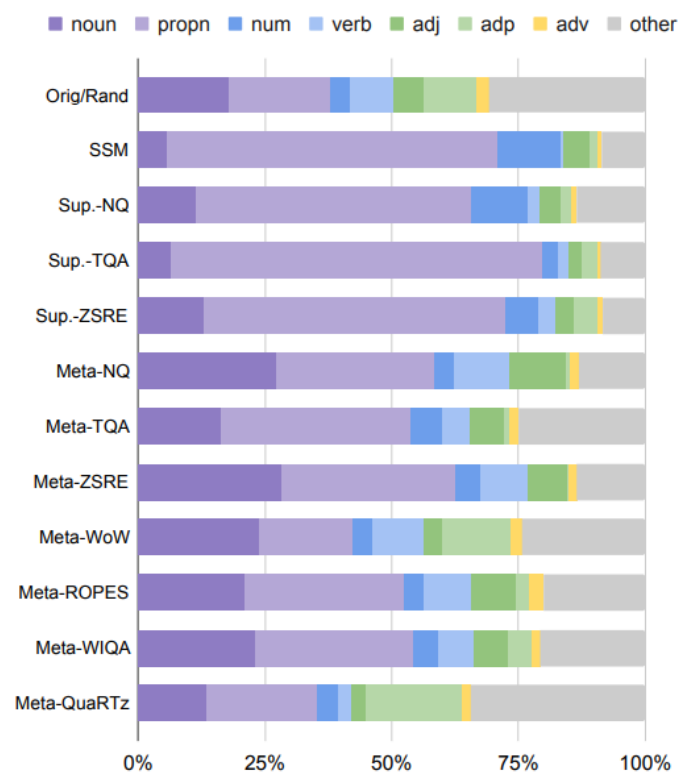


Please check out the full result tables in our paper!

	TQA	WQ	NQ
BART-Base	21.82 \pm 1.15	26.23 \pm 0.05	23.72 \pm 0.25
+Orig		AY2	ZSRE
+Rand		EM	EM
+SSM			WoW
+Supervised			F1
+Supervised		81.07	1.80
+Supervised			15.14
+Supervised		ROPES	WIQA
+Supervised			QuaRTz
+Supervised		BART-Base	46.60 \pm 0.48
+Supervised		+Orig	71.18 \pm 1.12
+Supervised		+SSM	62.80 \pm 1.16
+Supervised		+Rand	43.68 \pm 0.67
+Supervised		+SSM	73.06 \pm 0.72
+Supervised		+Meta-learned-ROPES	63.35 \pm 0.52
+Supervised		+Meta-learned-WIQA	44.59 \pm 1.15
+Supervised		+Meta-learned-QuaRTz	70.55 \pm 0.42
+Supervised		+Supervised	63.31 \pm 1.74
+Supervised		+Supervised	50.51 \pm 1.15
+Supervised		+Supervised	69.31 \pm 0.77
+Supervised		+Supervised	64.41 \pm 1.04
+Supervised		+Supervised	53.71 \pm 2.33
+Supervised		+Supervised	73.05 \pm 0.98
+Supervised		+Supervised	62.93 \pm 1.28
+Supervised		+Supervised	48.30 \pm 0.69
+Supervised		+Supervised	72.38 \pm 0.37
+Supervised		+Supervised	63.14 \pm 1.26
+Supervised		+Supervised	49.01 \pm 1.92
+Supervised		+Supervised	72.65 \pm 0.53
+Supervised		+Supervised	63.69 \pm 0.48
+Supervised		+Supervised	21.10 \pm 0.31
+Supervised		+Supervised	20.11 \pm 0.74
+Supervised		+Supervised	21.29 \pm 0.28

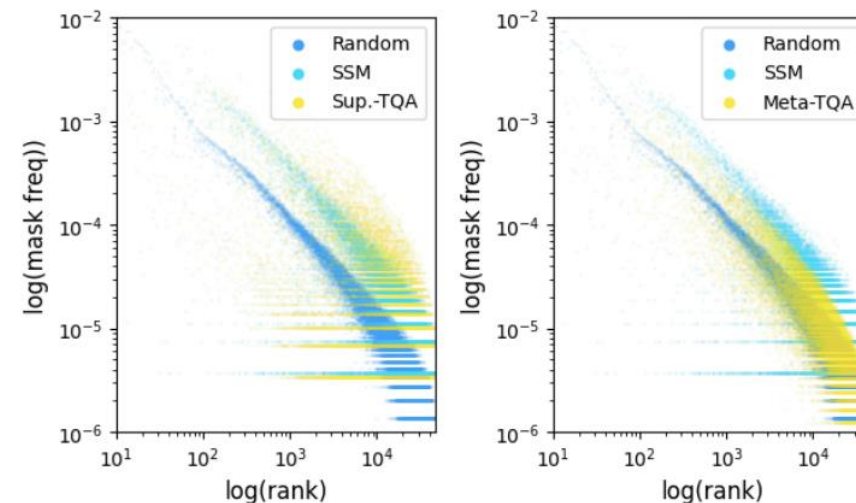
Analysis – Quantitative Analysis

Relation to Part-of-speech Tags



The learned masks are more *customized* towards the downstream tasks.

Relation to Token Frequency



(a) Supervised-TriviaQA (b) Meta-learned-TriviaQA

Random → approximates a Zipfian Distribution
SSM → Weaker preference for Zipfian
Learned Policies → Even weaker preference

Conclusions

- We introduce ***an analysis protocol*** to study the influence brought by different masking policies.
- We describe ***two methods (supervised/meta-learned masking policies)*** to learn masking policies.
- We conduct a large-scale ***empirical study*** with ***9 NLP tasks*** and ***3 categories of masking policies***.
- We ...
 - identify several ***successful cases*** of intermediate pre-training,
 - offer ***in-depth analysis and insights*** for the masking policies we used,
 - discuss the ***pros and cons*** of learned masking policies,
 - ***summarize several suggestions and tips*** for researchers who wish to adopt intermediate pre-training in their applications.