

# SELF-RAG: LEARNING TO RETRIEVE, GENERATE, AND CRITIQUE THROUGH SELF-REFLECTION

**Akari Asai<sup>†</sup>, Zeqiu Wu<sup>†</sup>, Yizhong Wang<sup>†§</sup>, Avirup Sil<sup>‡</sup>, Hannaneh Hajishirzi<sup>†§</sup>**

<sup>†</sup>University of Washington    <sup>§</sup>Allen Institute for AI    <sup>‡</sup>IBM Research AI

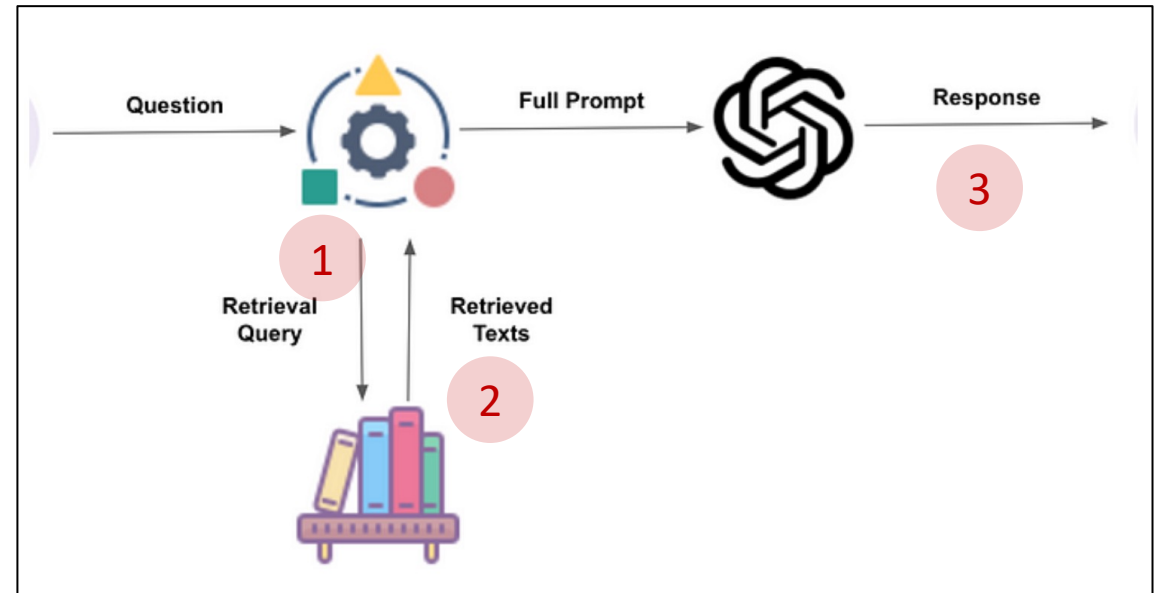
{akari, zeqiuwu, yizhongw, hannaneh}@cs.washington.edu, avi@us.ibm.com

ICLR2024 8886

# Motivation & Preliminary

- Why is Retrieval Augmented Generation (RAG) important?
  - Reduce factual errors
  - Improve helpfulness/usefulness
- Current Challenges of RAG
  1. May hinder the versatility of LLMs
    - **Learn when to retrieve**
  2. May introduce unnecessary or off-topic passages
    - **Learn what to retrieve (Part of critique)**
  3. (weak/small) LLMs may not know how to use the retrieved knowledge.
    - **Learn to generate (given retrieved data)**

## SELF-RAG: LEARNING TO RETRIEVE, GENERATE, AND CRITIQUE THROUGH SELF-REFLECTION



### Special Reflection Tokens to Learn

Type	Input	Output	Definitions
<b>Retrieve</b>	$x / x, y$	{yes, no, continue}	Decides when to retrieve with $\mathcal{R}$
<b>ISREL</b>	$x, d$	{ <b>relevant</b> , irrelevant}	$d$ provides useful information to solve $x$ .
<b>ISSUP</b>	$x, d, y$	{ <b>fully supported</b> , partially supported, no support}	All of the verification-worthy statement in $y$ is supported by $d$ .
<b>ISUSE</b>	$x, y$	{5, 4, 3, 2, 1}	$y$ is a useful response to $x$ .

x: query, y: output, d: documents

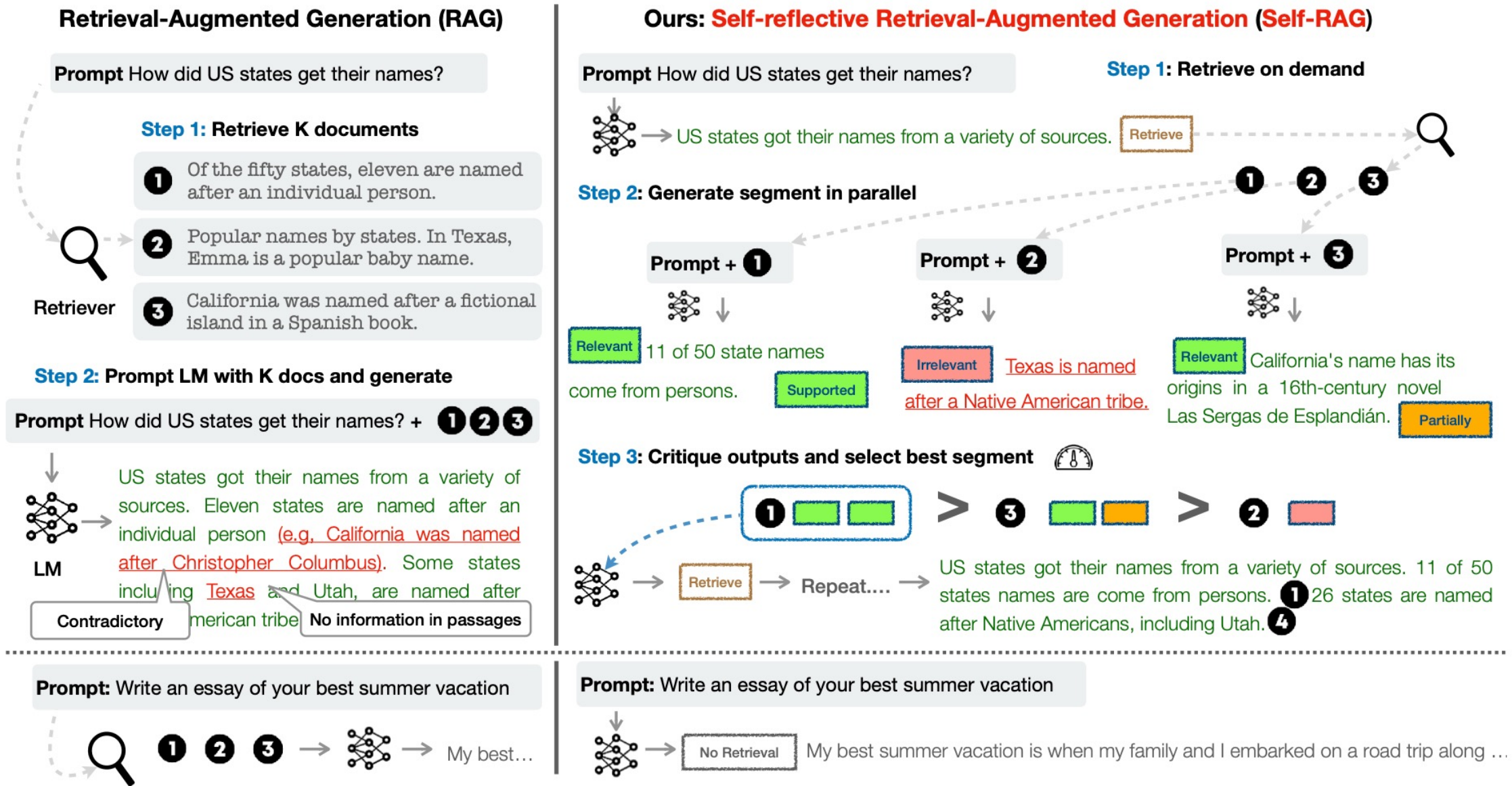


Figure 1: Overview of SELF-RAG. SELF-RAG learns to retrieve, critique, and generate text passages to enhance overall generation quality, factuality, and verifiability.

# Method: Training

## Learning to Retrieve\Critique

- Why?
  - Call GPT4 API is too expensive and hard to reproduce.
- Solution
  - Knowledge Distill from GPT4

$$\max_{\mathcal{C}} \mathbb{E}_{((x,y),r) \sim \mathcal{D}_{critic}} \log p_{\mathcal{C}}(r|x,y), \quad r \text{ for reflection tokens.}$$

Type	Input	Output	Definitions
<b>Retrieve</b>	$x / x, y$	{yes, no, continue}	Decides when to retrieve with $\mathcal{R}$
<b>ISREL</b>	$x, d$	{ <b>relevant</b> , irrelevant}	$d$ provides useful information to solve $x$ .
<b>ISSUP</b>	$x, d, y$	{ <b>fully supported</b> , partially supported, no support}	All of the verification-worthy statement in $y$ is supported by $d$ .
<b>ISUSE</b>	$x, y$	{5, 4, 3, 2, 1}	$y$ is a useful response to $x$ .

## Learning to Generate

- Use the learned Critic to augment output

- Then SFT  $\max_{\mathcal{M}} \mathbb{E}_{(x,y,r) \sim \mathcal{D}_{gen}} \log p_{\mathcal{M}}(y, r|x).$

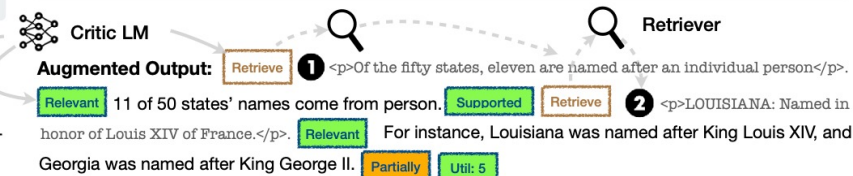
**Input:** Write an essay of your best summer vacation

**Output:** My best summer vacation was a magical escape to the coastal town of Santorini. The azure waters, charming white-washed building are unforgettable.

**Augmented Output:** No Retrieval My best summer vacation was a magical escape to the coastal town of Santorini. No Retrieval The azure waters, charming white-washed building are unforgettable experience. Util: 5

**Input:** How did US states get their names?

**Output:** 1 of 50 states names come from persons. For instance, Louisiana was named in honor of King Louis XIV of France and Georgia was named after King George II.



Dataset name	category	Data source	the number of instances
GPT-4 Alpaca	Instruction-following	Open-Instruct	26,168
Stanford Alpaca	Instruction-following	Open-Instruct	25,153
FLAN-V2	Instruction-following	Open-Instruct	17,817
ShareGPT	Instruction-following	Open-Instruct	13,406
Open Assistant 1	Instruction-following	Open-Instruct	9,464
Wizard of Wikipedia	Knowledge-intensive	KILT	17,367
Natural Questions	Knowledge-intensive	KILT	15,535
FEVER	Knowledge-intensive	KILT	9,966
OpenBoookQA	Knowledge-intensive	HF Dataset	4,699
Arc-Easy	Knowledge-intensive	HF Dataset	2,147
ASQA	Knowledge-intensive	ASQA	3,897

Total: 150k

### Algorithm 2 SELF-RAG Training

- 1: **Input** input-output data  $\mathcal{D} = \{X, Y\}$ , generator  $\mathcal{M}, \mathcal{C} \theta$
- 2: Initialize  $\mathcal{C}$  with a pre-trained LM
- 3: Sample data  $\{X^{sample}, Y^{sample}\} \sim \{X, Y\}$  ▷ **Training Critic LM (Section 3.2.1)**
- 4: **for**  $(x, y) \in (X^{sample}, Y^{sample})$  **do** ▷ Data collections for  $\mathcal{C}$
- 5:     Prompt GPT-4 to collect a reflection token  $r$  for  $(x, y)$
- 6:     Add  $\{(x, y, r)\}$  to  $\mathcal{D}_{critic}$
- 7: Update  $\mathcal{C}$  with next token prediction loss ▷ Critic learning; Eq. 1
- 8: Initialize  $\mathcal{M}$  with a pre-trained LM ▷ **Training Generator LM (Section 3.2.2)**
- 9: **for**  $(x, y) \in (X, Y)$  **do** ▷ Data collection for  $\mathcal{M}$  with  $\mathcal{D}_{critic}$
- 10:     Run  $\mathcal{C}$  to predict  $r$  given  $(x, y)$
- 11:     Add  $(x, y, r)$  to  $\mathcal{D}_{gen}$
- 12: Update  $\mathcal{M}$  on  $\mathcal{D}_{gen}$  with next token prediction loss ▷ Generator LM learning; Eq. 2

# Method: Inference

- Step1: Adaptive retrieval with threshold.
- Step2: Call retriever  $R$
- Step3: Tree-decoding with critique tokens.
  - Beam-Search

$$f(y_t, d, \text{Critique}) = p(y_t|x, d, y_{<t}) + \mathcal{S}(\text{Critique}), \text{ where} \quad (3)$$

$$\mathcal{S}(\text{Critique}) = \sum_{G \in \mathcal{G}} w^G s_t^G \text{ for } \mathcal{G} = \{\text{ISREL}, \text{ISUP}, \text{ISUSE}\}, \quad (4)$$

where  $s_t^G = \frac{p_t(\hat{r})}{\sum_{i=1}^{N^G} p_t(r_i)}$  stands for the generation probability of the most desirable reflection token  $\hat{r}$  (e.g.,  $\text{ISREL} = \text{Relevant}$ ) for the critique token type  $G$  with  $N^G$  distinct tokens (that represent

Sensitive to  $w$ , a customizable feature

## Algorithm 1 SELF-RAG Inference

**Require:** Generator LM  $\mathcal{M}$ , Retriever  $\mathcal{R}$ , Large-scale passage collections  $\{d_1, \dots, d_N\}$

- 1: **Input:** input prompt  $x$  and preceding generation  $y_{<t}$ , **Output:** next output segment  $y_t$
- 2:  $\mathcal{M}$  predicts **Retrieve** given  $(x, y_{<t})$
- 3: **if** **Retrieve** == Yes **then**
- 4: Retrieve relevant text passages  $\mathbf{D}$  using  $\mathcal{R}$  given  $(x, y_{t-1})$  ▷ **Retrieve**
- 5:  $\mathcal{M}$  predicts **ISREL** given  $x, d$  and  $y_t$  given  $x, d, y_{<t}$  for each  $d \in \mathbf{D}$  ▷ **Generate**
- 6:  $\mathcal{M}$  predicts **ISUP** and **ISUSE** given  $x, y_t, d$  for each  $d \in \mathbf{D}$  ▷ **Critique**
- 7: Rank  $y_t$  based on **ISREL**, **ISUP**, **ISUSE** ▷ Detailed in Section 3.3
- 8: **else if** **Retrieve** == No **then**
- 9:  $\mathcal{M}_{gen}$  predicts  $y_t$  given  $x$  ▷ **Generate**
- 10:  $\mathcal{M}_{gen}$  predicts **ISUSE** given  $x, y_t$  ▷ **Critique**

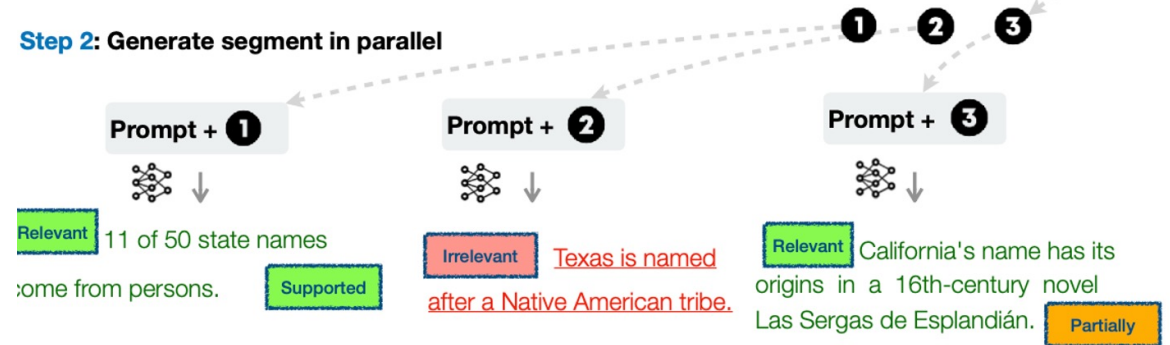
## Ours: Self-reflective Retrieval-Augmented Generation (Self-RAG)

Prompt How did US states get their names?

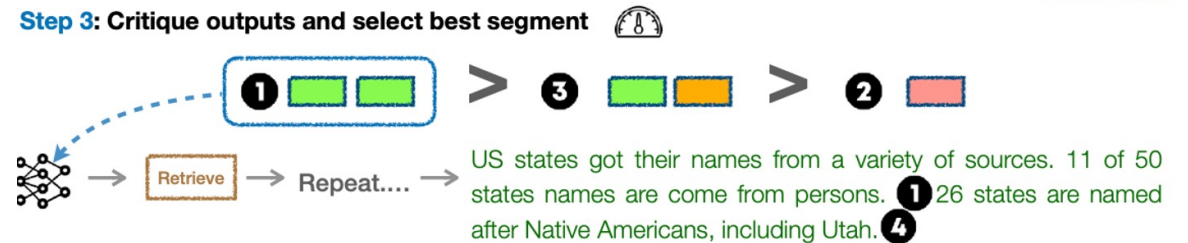
Step 1: Retrieve on demand



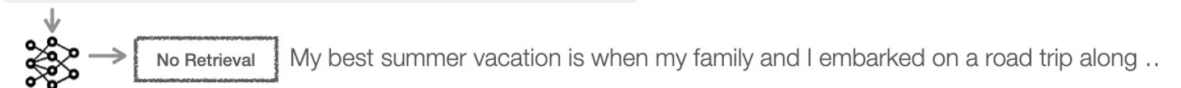
Step 2: Generate segment in parallel



Step 3: Critique outputs and select best segment



Prompt: Write an essay of your best summer vacation



# Experiments: Main results

Setting:  
 Critic C: 7b-base  
 Generator M: 7b and 13b

LM	Short-form		Closed-set		Long-form generations (with citations)					
	PopQA (acc)	TQA (acc)	Pub (acc)	ARC (acc)	Bio (FS)	(em)	(rg)	ASQA (mau)	(pre)	(rec)
<i>LMs with proprietary data</i>										
Llama2-c <sub>13B</sub>	20.0	59.3	49.4	38.4	55.9	22.4	29.6	28.6	–	–
Ret-Llama2-c <sub>13B</sub>	51.8	59.8	52.1	37.9	79.9	32.8	34.8	43.8	19.8	36.1
ChatGPT	29.3	<b>74.3</b>	70.1	<b>75.3</b>	71.8	35.3	36.2	68.8	–	–
Ret-ChatGPT	50.8	65.7	54.7	<b>75.3</b>	–	<b>40.7</b>	<b>39.9</b>	<b>79.7</b>	65.1	<b>76.6</b>
Perplexity.ai	–	–	–	–	71.2	–	–	–	–	–
<i>Baselines without retrieval</i>										
Llama2 <sub>7B</sub>	14.7	30.5	34.2	21.8	44.5	7.9	15.3	19.0	–	–
Alpaca <sub>7B</sub>	23.6	54.5	49.8	45.0	45.8	18.8	29.4	61.7	–	–
Llama2 <sub>13B</sub>	14.7	38.5	29.4	29.4	53.4	7.2	12.4	16.0	–	–
Alpaca <sub>13B</sub>	24.4	61.3	55.5	54.9	50.2	22.9	32.0	70.6	–	–
CoVE <sub>65B</sub> *	–	–	–	–	71.2	–	–	–	–	–
<i>Baselines with retrieval</i>										
Toolformer* <sub>6B</sub>	–	48.8	–	–	–	–	–	–	–	–
Llama2 <sub>7B</sub>	38.2	42.5	30.0	48.0	78.0	15.2	22.1	32.0	2.9	4.0
Alpaca <sub>7B</sub>	46.7	64.1	40.2	48.0	76.6	30.9	33.3	57.9	5.5	7.2
Llama2-FT <sub>7B</sub>	48.7	57.3	64.3	65.8	78.2	31.0	35.8	51.2	5.0	7.5
SAIL* <sub>7B</sub>	–	–	69.2	48.4	–	–	–	–	–	–
Llama2 <sub>13B</sub>	45.7	47.0	30.2	26.0	77.5	16.3	20.5	24.7	2.3	3.6
Alpaca <sub>13B</sub>	46.1	66.9	51.1	57.6	77.7	<b>34.8</b>	36.7	56.6	2.0	3.8
<b>Our SELF-RAG</b> <sub>7B</sub>	54.9	66.4	72.4	67.3	<b>81.2</b>	30.0	35.7	<b>74.3</b>	66.9	67.8
<b>Our SELF-RAG</b> <sub>13B</sub>	<b>55.8</b>	<b>69.3</b>	<b>74.5</b>	<b>73.1</b>	80.2	31.7	<b>37.0</b>	71.6	<b>70.3</b>	<b>71.3</b>

# Experiments: Ablations

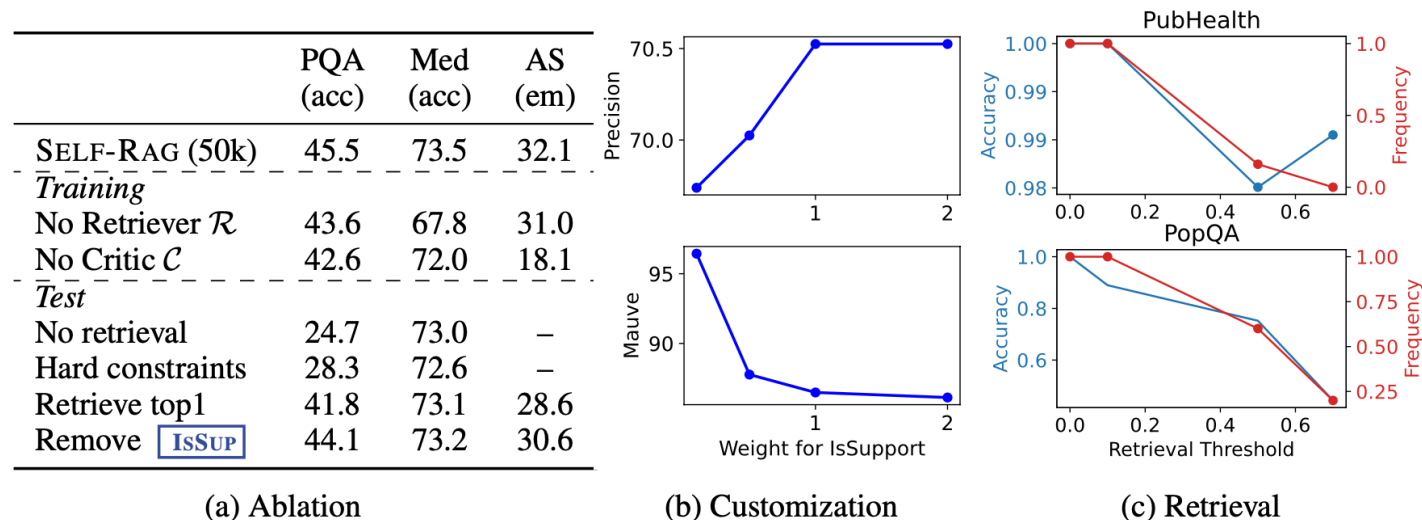


Figure 3: **Analysis on SELF-RAG:** (a) **Ablation studies** for key components of SELF-RAG training and inference based on our 7B model. (b) **Effects of soft weights** on ASQA citation precision and Mauve (fluency). (c) **Retrieval frequency** and *normalized* accuracy on PubHealth and PopQA.

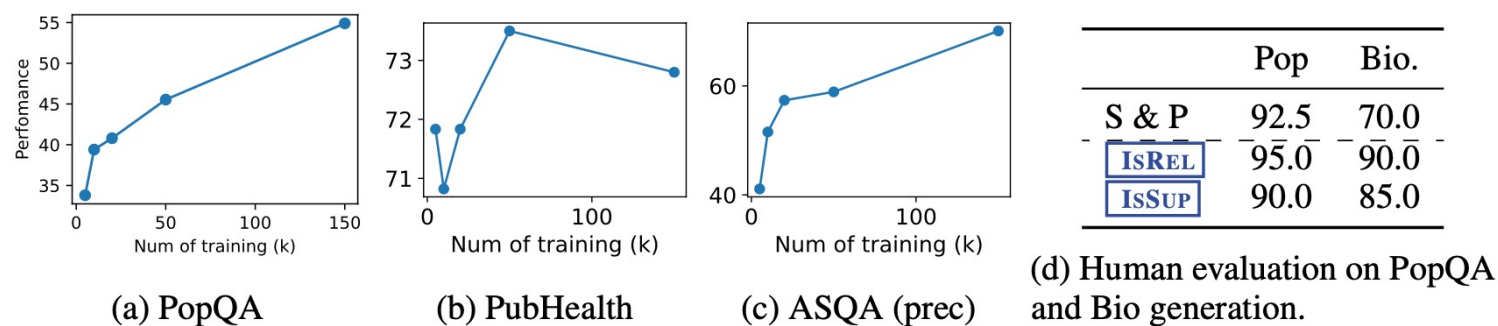


Figure 4: **Training scale and Human analysis:** (a) (b) (c) **Training scale analysis** shows the effect of the training data scale on PopQA, PubHealth and ASQA (citation precision), respectively. (d) **Human analysis** on SELF-RAG outputs as well as reflection tokens.