



# Training Large Language Models for Retrieval-Augmented Question Answering through Backtracking Correction

---

**Huawen Feng**, Zekun Yao, Junhao Zheng, Qianli Ma

School of Computer Science and Engineering,

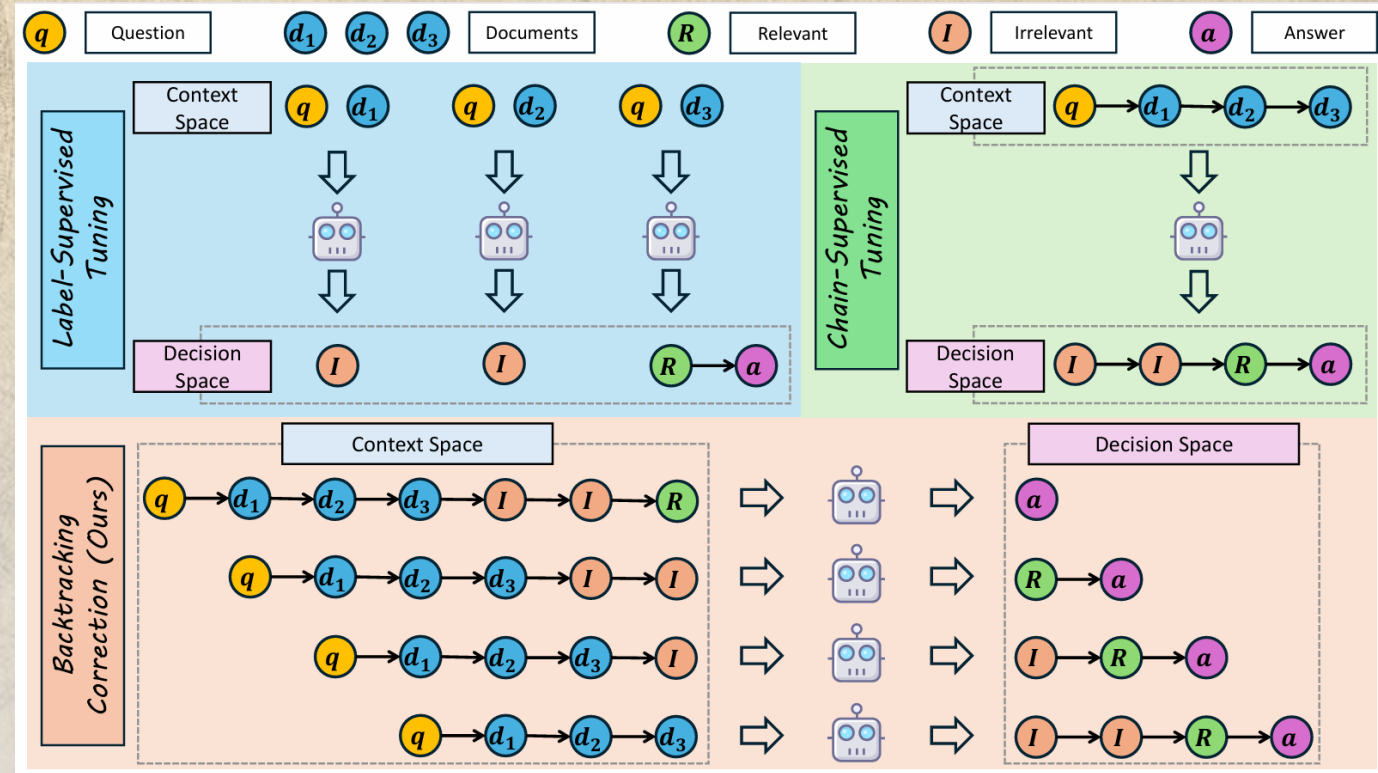
South China University of Technology,

Guangzhou, Guangdong, China



# BACKGROUND

- Despite recent progress in Retrieval-Augmented Generation (RAG) achieved by large language models (LLMs), retrievers often recall uncorrelated documents, regarded as "noise" during subsequent text generation. To address this, current methods train LLMs to distinguish between relevant and irrelevant documents using labeled data, enabling them to select the most likely relevant ones as context.
  - Label-Supervised Tuning**
  - Chain-Supervised Tuning**

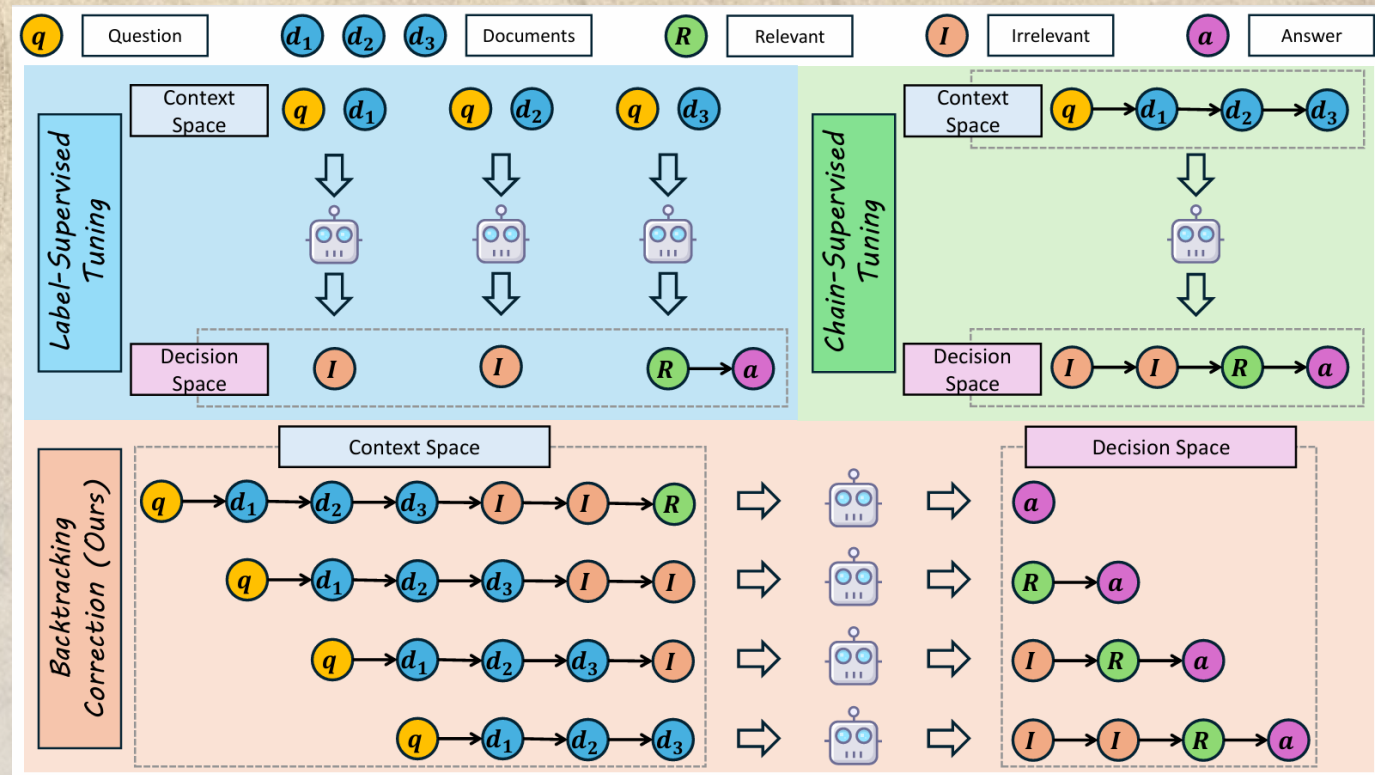




# LIMITS



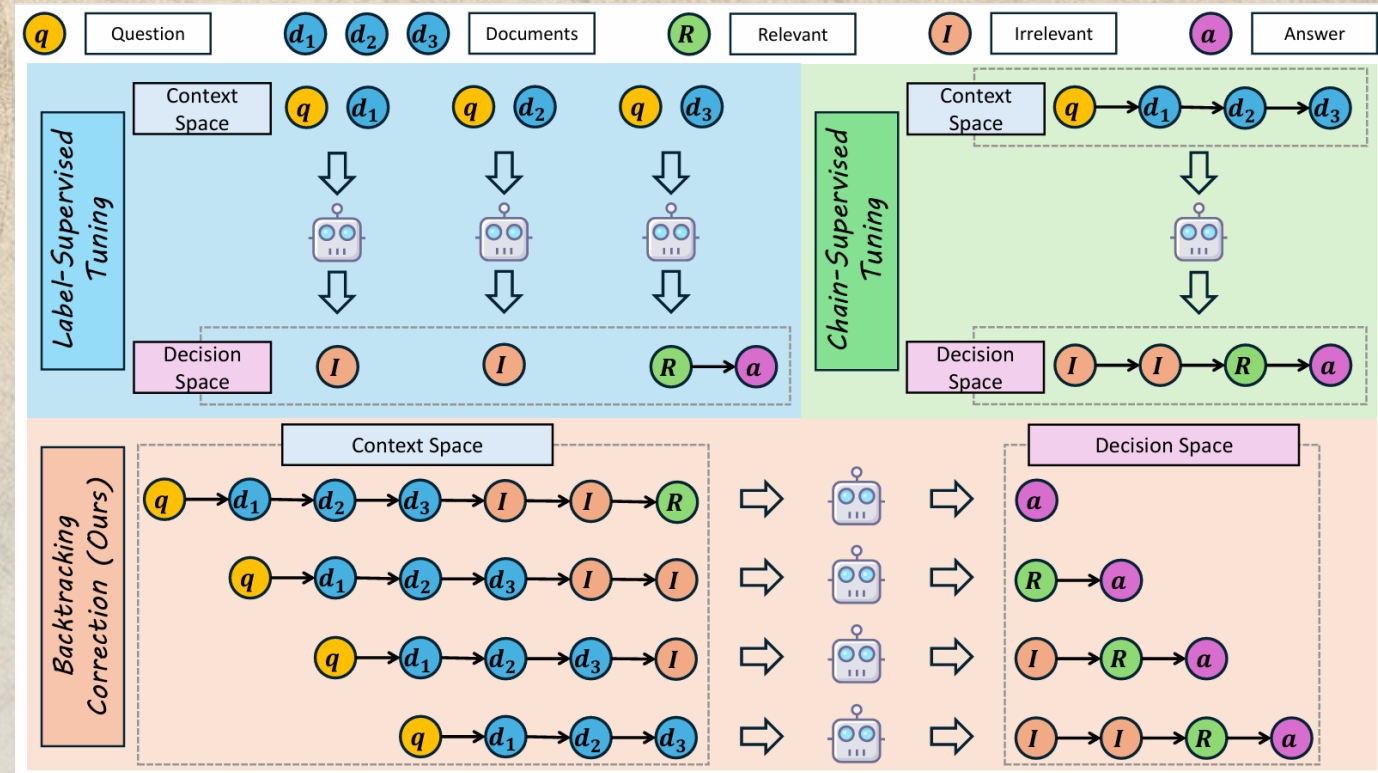
- **Label-Supervised Tuning**—  
The sparse supervision in the decision-making process limits the model's ability to understand why a retrieved document is relevant or irrelevant.
- **Chain-Supervised Tuning**—  
Stepwise-annotated data is challenging to collect, and much of it depends on proprietary language models (e.g., ChatGPT, GPT-4).
- Both of them requires **extensive training data with high-quality annotations**, significantly increasing costs.





# MOTIVATION

We apply **stepwise preference optimization** to reasoning chains. Starting from the final state, we sample model-generated errors and employ self-correction to optimize the current state through Reinforcement Learning (RL), then backtrack to the previous state. Unlike Chain-Supervised Tuning which learns the entire reasoning path, our method optimizes each step by focusing on decision-making based on the current state, as the remaining steps have already been learned. This approach gradually reduces the context space while expanding the decision space, following an easy-to-hard progression.









# DERIVATION

Player A:

$$\arg \max_{\theta} \mathbb{E}_{s_{t-1}^+, o_t \sim \Pi_{\theta}(s_{t-1}^+), o_t^+ \sim \Pi_{\theta}(s_{t-1}^+, o_t, L)} [r(s_{t-1}^+, o_t^+) - r(s_{t-1}^+, o_t)]$$

$$\arg \min_{\theta} -\mathbb{E}_{s_{t-1}^+, o_t \sim \Pi_{\theta}(s_{t-1}^+), o_t^+ \sim \Pi_{\theta}(s_{t-1}^+, o_t, L)} [\log \sigma(r(s_{t-1}^+, o_t^+) - r(s_{t-1}^+, o_t))]$$

Player B:

$$\begin{aligned} \arg \max_{\theta} \mathbb{E}_{s_{t-1}^+, o_t \sim \Pi_{\theta}(s_{t-1}^+), o_t^+ \sim \Pi_{\theta}(s_{t-1}^+, o_t, L)} [r(s_{t-1}^+, o_t^+)] \\ - \beta \mathbb{D}_{KL} [\Pi_{\theta}(o_t^+ | s_{t-1}^+) || \Pi_{ref}(o_t^+ | s_{t-1}^+)] \end{aligned}$$

$$\Pi_{\theta}(o_t^+ | s_{t-1}^+) = \frac{1}{S(s_{t-1}^+)} \Pi_{ref}(o_t^+ | s_{t-1}^+) \exp\left(\frac{1}{\beta} r(s_{t-1}^+, o_t^+)\right)$$

$$\Pi_{\theta}(o_t | s_{t-1}^+) = \frac{1}{S(s_{t-1}^+)} \Pi_{ref}(o_t | s_{t-1}^+) \exp\left(\frac{1}{\beta} r(s_{t-1}^+, o_t)\right)$$





# DERIVATION

---

$$r(s_{t-1}^+, o_t^+) = \beta \log \frac{\Pi_{\theta}(o_t^+ | s_{t-1}^+)}{\Pi_{ref}(o_t^+ | s_{t-1}^+)} \quad r(s_{t-1}^+, o_t) = \beta \log \frac{\Pi_{\theta}(o_t | s_{t-1}^+)}{\Pi_{ref}(o_t | s_{t-1}^+)}$$

$$\mathcal{L}_{BC} = \arg \min_{\theta} -\mathbb{E} \left[ \log \sigma \left( \beta \log \frac{\Pi_{\theta}(o_t^+ | s_{t-1}^+)}{\Pi_{ref}(o_t^+ | s_{t-1}^+)} - \beta \log \frac{\Pi_{\theta}(o_t | s_{t-1}^+)}{\Pi_{ref}(o_t | s_{t-1}^+)} \right) \right]$$

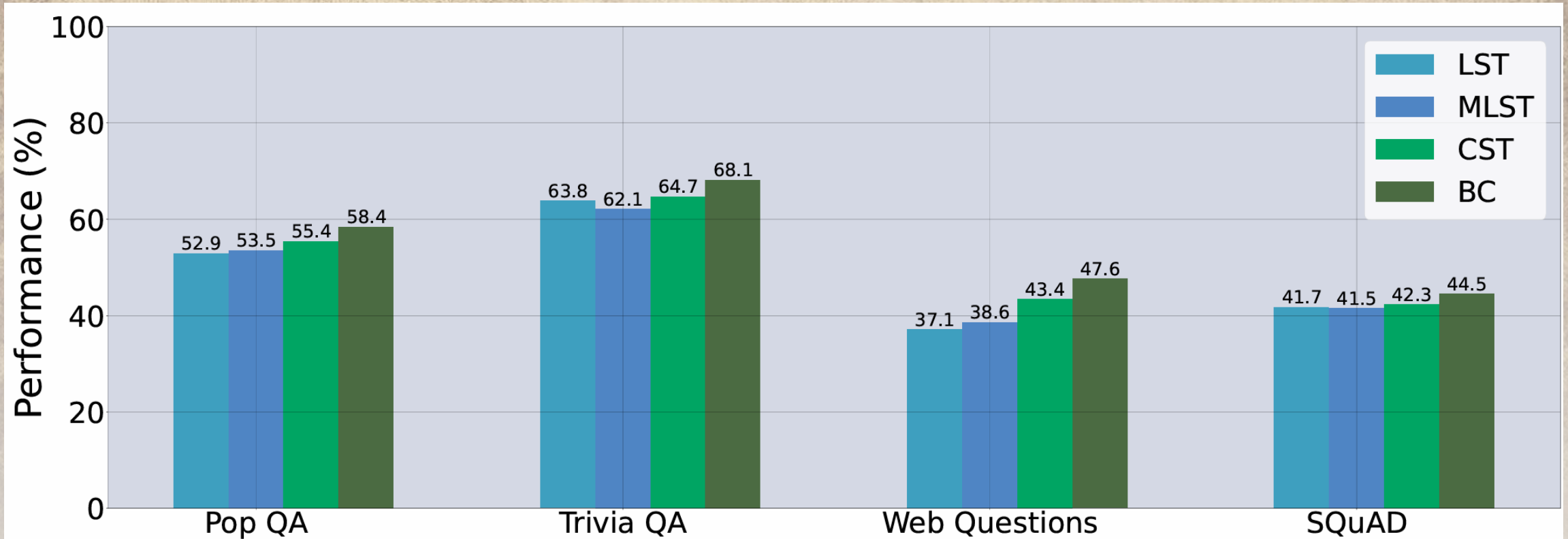


# EXPERIMENT

		Pop QA	Trivia QA	Web Questions	SQuAD
Proprietary	ChatGPT	24.7	78.2	57.0	28.5
	GPT-4	38.7	83.4	61.6	35.1
Base	LLaMA2-7B-base	4.5	10.9	6.1	3.7
	LLaMA3-8B-base	5.7	10.8	4.4	3.8
Fine-tuned w/o Retrieval	ChatGLM3-6B	9.5	19.6	14.9	4.3
	Mistral-7B	25.2	66.1	56.8	25.9
	BaiChuan2-7B-chat	25.7	40.7	38.6	13.1
	LLaMA2-7B-chat	25.1	58.7	48.6	19.1
	Alpaca-7B	25.6	49.4	39.6	14.1
Fine-tuned w/ Retrieval	ChatGLM3-6B	42.6	35.7	24.3	19.5
	ToolFormer-6B	-	48.8	26.3	33.8
	Mistral-7B	59.1	71.2	57.6	42.0
	BaiChuan2-7B-chat	56.2	63.3	50.0	38.7
	LLaMA2-7B-chat	54.0	63.1	45.3	37.2
	Alpaca-7B	54.2	58.6	47.3	32.6
	Self-RAG-7B	53.8	62.6	32.4	26.5
	REAR	52.9	71.6	38.4	43.8
	<i>RobustLM-13B</i>	<i>49.1</i>	<i>62.0</i>	<i>27.3</i>	<i>27.4</i>
Backtracking	+LLaMA2-7B-base	58.4	68.1	47.6	44.5
Correction	+LLaMA3-8B-base	59.3	70.1	49.3	44.9

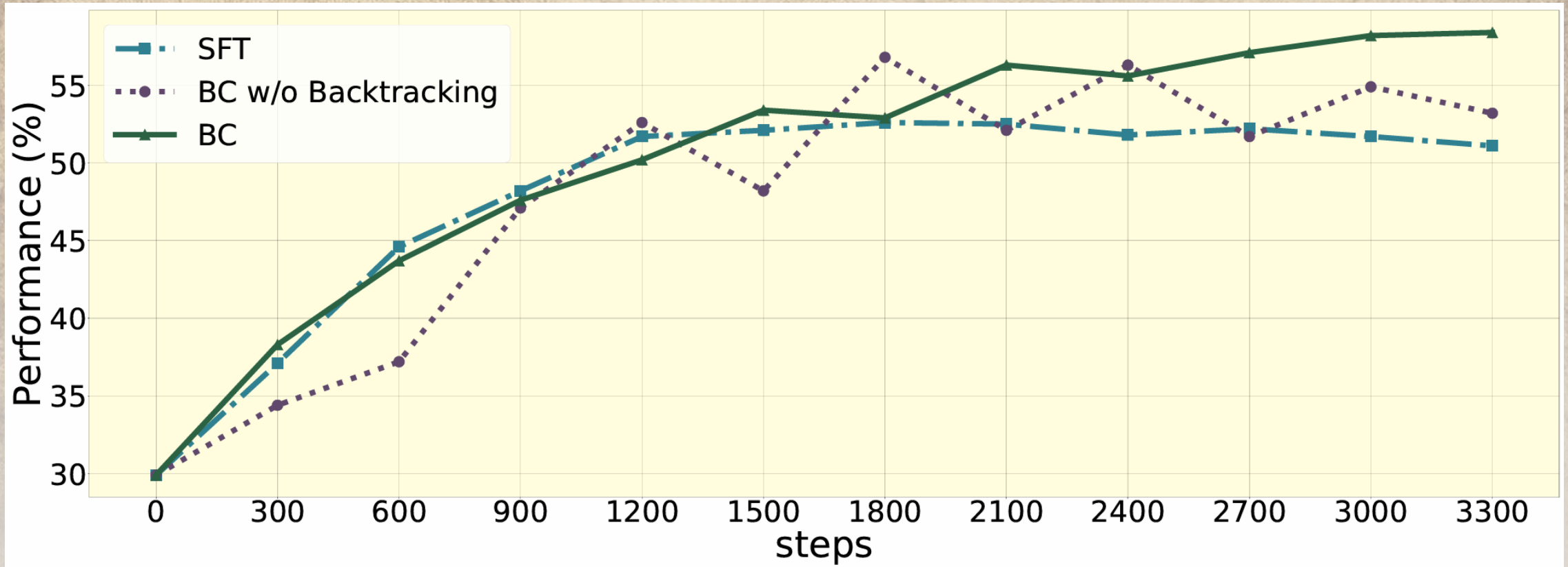


# EXPERIMENT





# EXPERIMENT





The background of the slide features a light beige, textured paper-like surface. On the left side, there is a faint, detailed illustration of a hot air balloon with a patterned envelope and a basket. On the right side, there is a faint illustration of a blimp or rigid airship. A solid red horizontal line is positioned above the main text.

# Thanks for your listening!