

# RoToR: Towards More Reliable Responses for Order-Invariant Inputs



Soyoung Yoon<sup>1\*</sup>, Dongha Ahn<sup>12</sup>, Youngwon Lee<sup>1</sup>, Minkyu Jung<sup>2</sup>, HyungJoo Jang<sup>2</sup>, Seung-won Hwang<sup>1†</sup>

<sup>1</sup>Seoul National University <sup>2</sup>Channel Corporation

<sup>†</sup>Corresponding Author <sup>\*</sup>Work done during an internship at Channel Corporation

Paper

Code

## Overview

We introduce **RoToR**: a **zero-shot** order-invariant Language Model, with (1) Global Sorting + Circular Position IDs and (2) Selective Routing for Mixed Inputs, which achieve SOTA robustness on 3 benchmarks & 25-45% lower FLOPs v.s. Baselines (PINE)

## Motivation: Positional Bias for Listwise Inputs

- Lost-in-the Middle (RAG)
- First-choice bias (75%) in LLM-as-a-judge
- MMLU rank shifts by 8 with shuffle -> need neutral handling for sets, tables, multiple-choice questions
- Zero-shot invariant LMs have been proposed as a solution, but had 2 limitations:

Which one is red?

A. Apple

B. Orange

C. Grape

Answer: ○

-> A. Apple

Which one is red?

A. Orange

B. Apple

C. Grape

Answer: ✗

-> A. Orange

## Limitations of previous Zero-shot order-invariant LMs

### - Limitation 1: Training and inference distribution mismatch

- PCW, Set-based prompting: No cross-segment content
- PINE: per-query sort ->  $O(O(n^2) + \text{instability})$
- Frequent ID changes cause **OOD behavior** -> drops its ability

### - Limitation 2: Fail to extend to real-life scenarios (order-invariant + order-sensitive)

- Does not consider, cannot be applied to hybrid cases (e.g., MMLU)

**PINE**: Bidirectional processing with Q-K similarity

- Has to obtain the same attention representation, regardless of initial ordering of segments
- Places query IDs last, sorts other segments in a order-invariant way

Challenges of PINE

- Frequent alterations on position IDs (layer, head, suffix, generation tokens)
- Computationally expensive
- Numerical Instability (arising from attention assignment)

Self-attention patterns (x = query, y = key) across order-invariant models

Example of order-sensitive and order-invariant cases

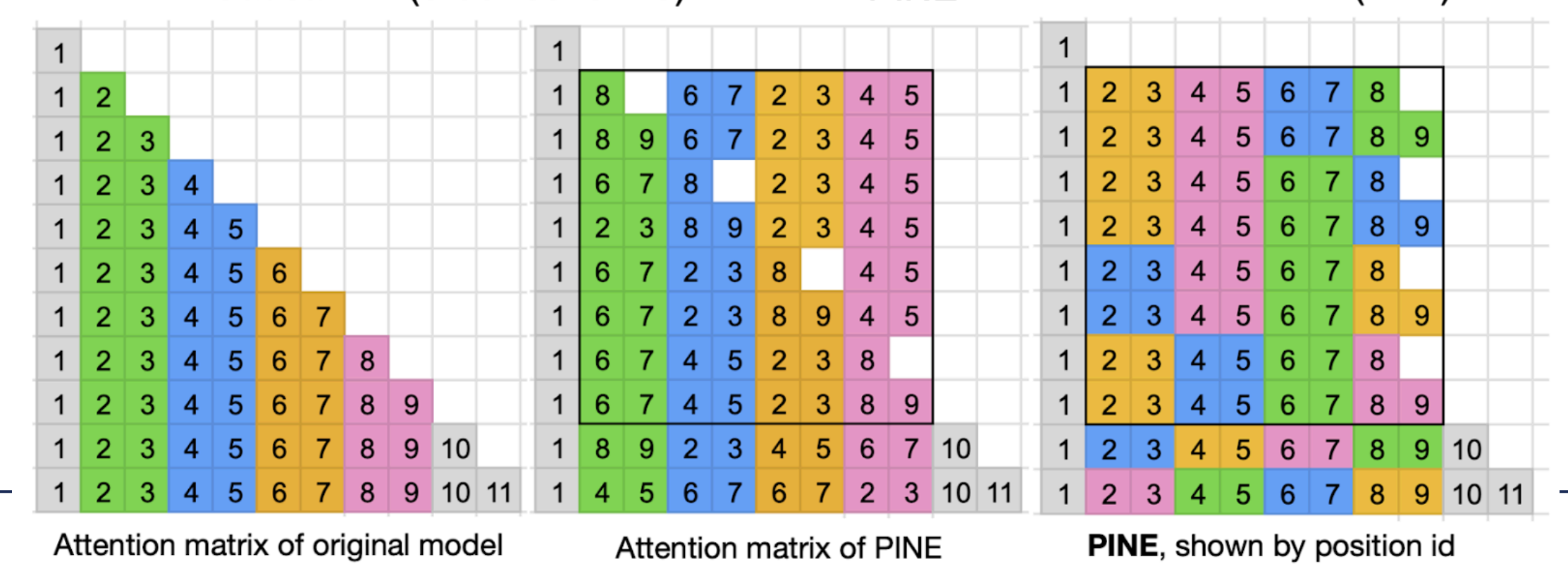
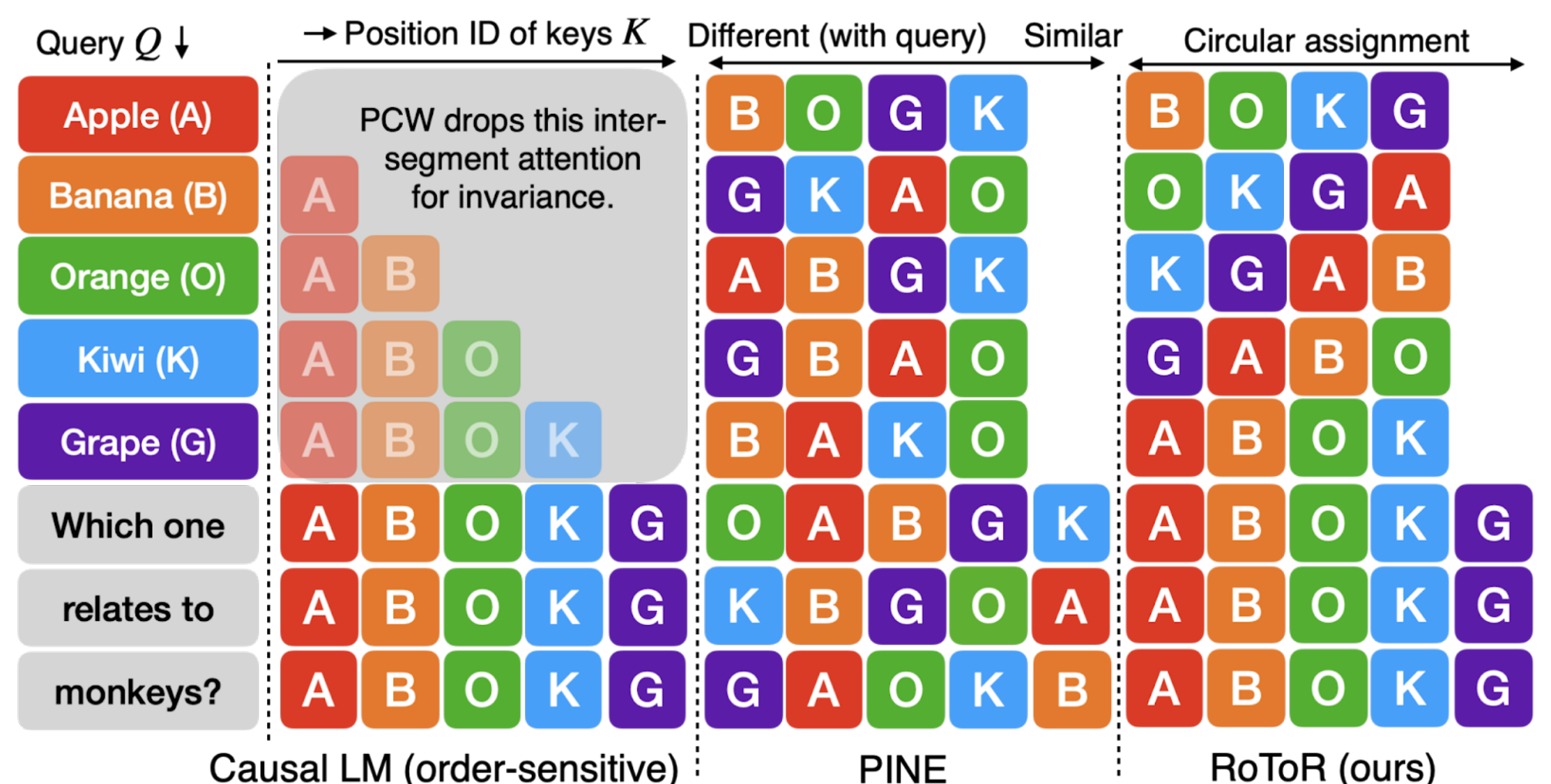
In 8085 name/names of the 16 bit registers is/are:

Order sensitive

A. stack pointer  
B. program counter  
C. both A and B.  
D. none of these

Order invariant

A. stack pointer  
B. program counter  
C. accumulator  
D. microprocessor



## Proposed Method: RoToR

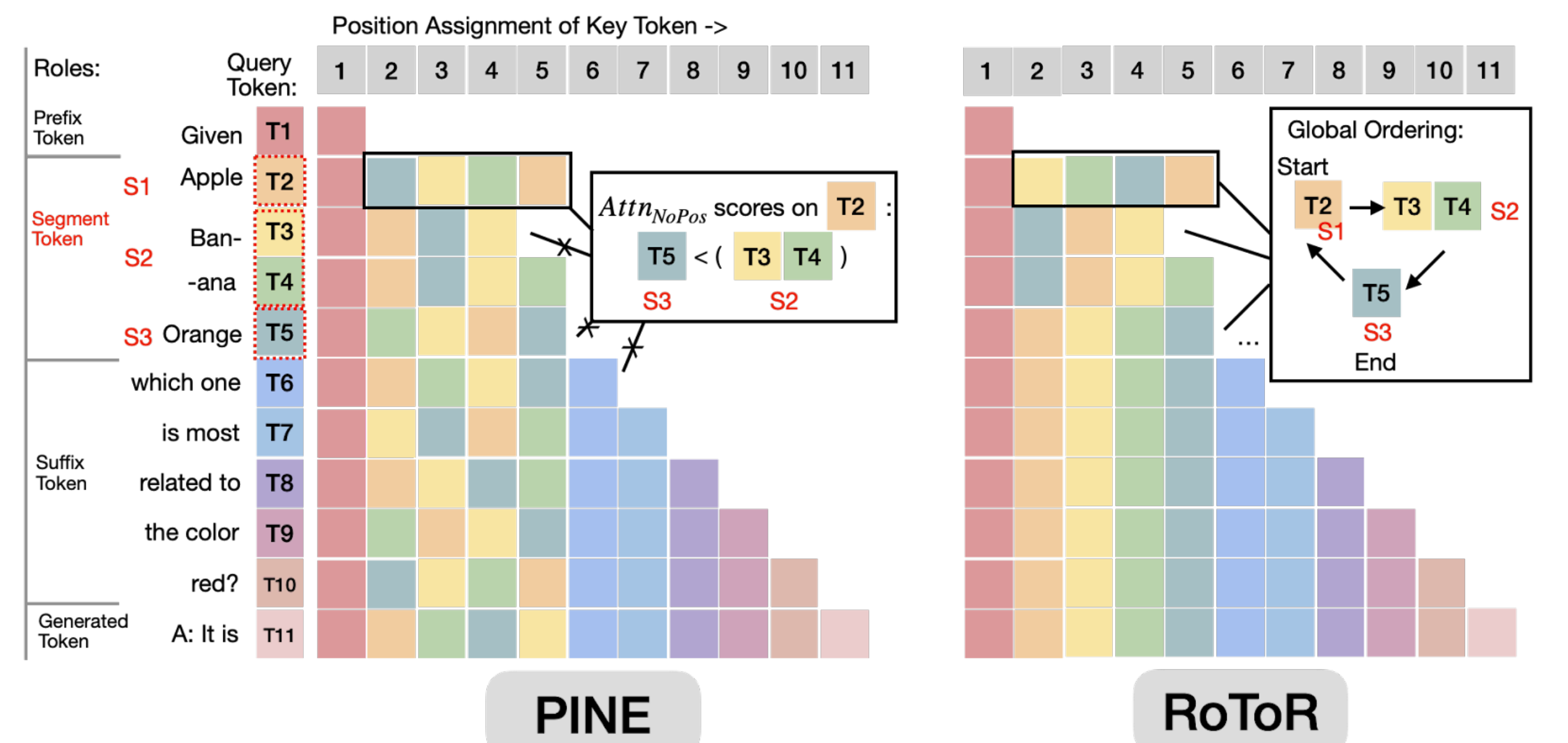
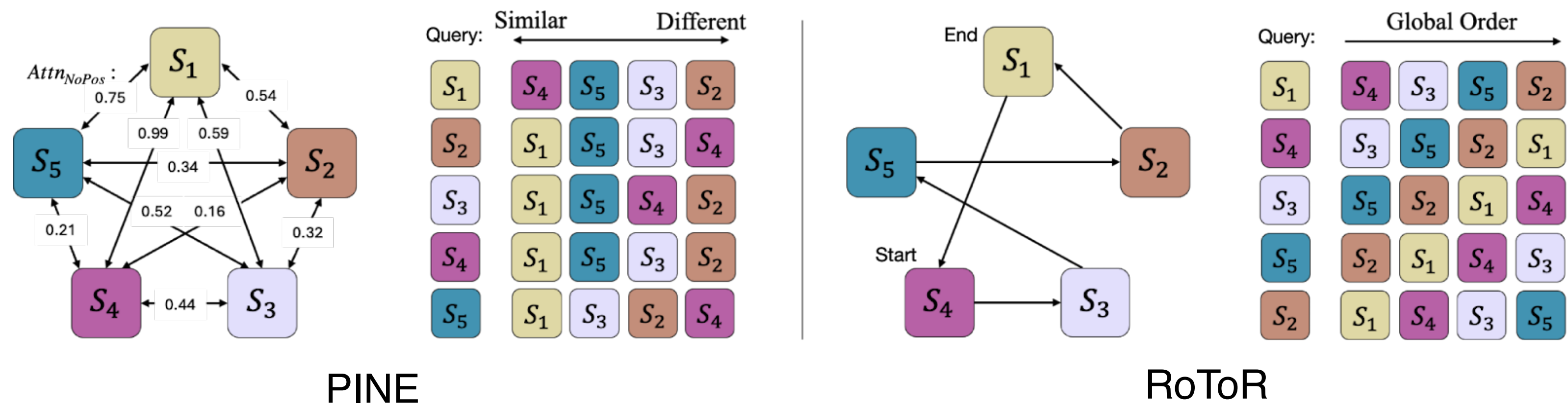
- Keep the bidirectional structure, but alter the position assignment in a simple and stable way!
- Define a single global ordering + circular assignment
- We propose three different global ordering methodology: (1) Lexical: Hierarchical sort on tokenized ID numbers, (2) Reranking (monot5), (3) Frequency-normalized

### 1. Training and inference distribution mismatch

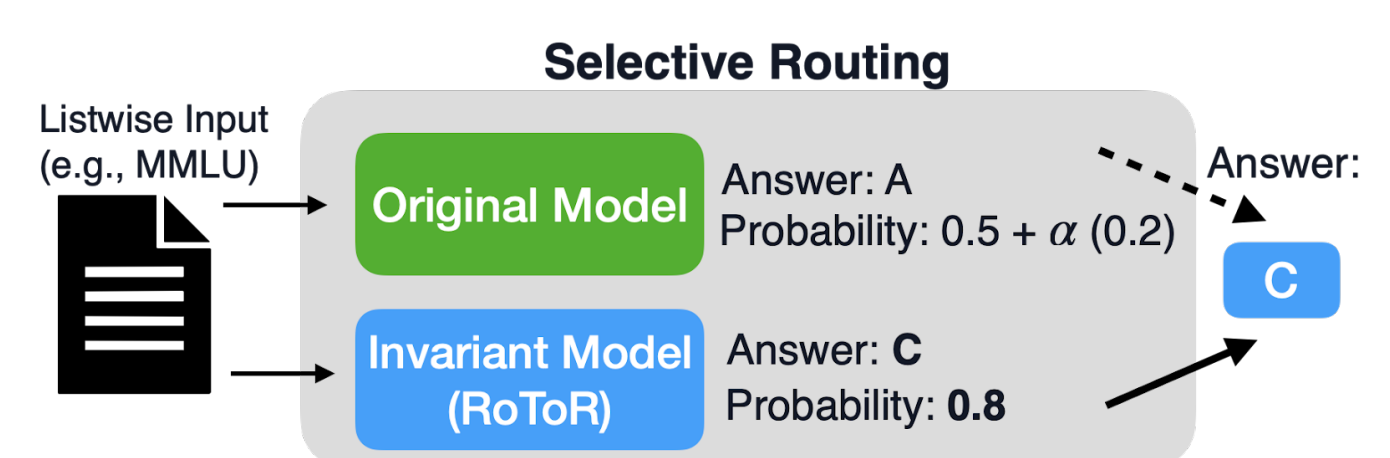
- Stable, order-invariant solution (RoToR)
- Query-agnostic global ordering with minimal positional ID modifications

### 2. Fail to extend to hybrid cases

- Selective Routing, which switches between original / invariant LMs based on confidence



**PINE**: query-dependent grid, **RoToR**: fixed order, rotate per query  
-> Stable IDs, zero collisions, less computation



## Experimental Setup

### - Benchmarks:

- Lost-in-the-Middle (LitM)
- Knowledge Graph QA (KGQA): Mintaka
- MMLU: selective routing cases
- LongBench: long context scenarios (Appendix)

### - Model backbones:

- Llama-3.1-8B/70B
- Qwen-1.5-4/7/72B-Chat

- **Metrics**: best\_subspan\_em (LitM), EM, F1, Acc. (KGQA), Acc. (MMLU)

- **Methods**: Original (order-sensitive), PCW, Set-based prompting, PINE, RoToR

## Efficiency

- Less computation:
  - Overhead FLOPs ↓ 98 % (72B)
- Faster:
  - E2E Latency ↓ 23-43 % on LitM
- Reduces OOD:
  - Perplexity ↓; collision rate 0 %

## KGQA

	Llama-3.1-8B-Instr.			Llama-3.1-70B-Instr.			Qwen1.5-4B-Chat			Qwen1.5-7B-Chat			Qwen1.5-72B-Chat		
Method	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1
<b>N = 30</b>															
<b>Initial, no shuffling of segments</b>															
Original	50.2	44.0	51.9	61.6	57.7	63.6	30.7	27.9	34.9	31.5	27.8	35.4	41.4	37.7	43.7
PINE	51.5	45.0	52.6	63.1	58.7	64.8	31.6	28.7	35.6	32.3	28.8	36.4	46.7	42.9	49.0
RoToR	<b>53.1</b>	<b>46.5</b>	<b>54.1</b>	<b>63.6</b>	<b>59.1</b>	<b>65.2</b>	32.0	29.0	35.7	<b>34.3</b>	<b>29.8</b>	<b>37.7</b>	<b>47.5</b>	<b>43.2</b>	<b>49.2</b>
RoToR-MonoT5	51.6	45.0	52.5	-	-	-	32.3	29.1	36.2	32.9	28.4	36.3	-	-	-
RoToR-Freq.	52.6	46.1	53.7	-	-	-	<b>32.3</b>	<b>29.2</b>	36.0	33.7	29.5	37.2	-	-	-
<b>After shuffling segments, averaged over 3 seeds</b>															
Original	49.5	43.3	51.0	62.1	57.8	64.0	30.1	27.5	34.7	31.4	27.3	35.0	41.0	37.6	43.6
→ stdev. (±)	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.41/0.34/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.37/0.40/0.33	0.37/0.40/0.33	0.37/0.40/0.33
PINE	51.8	45.2	52.8	63.3	58.8	64.9	31.5	28.7	35.6	32.3	28.8	35.7	46.9	<b>43.3</b>	<b>49.2</b>
→ stdev. (±)	0.05/0.07/0.16	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.18/0.20/0.20	0.18/0.20/0.20	0.18/0.20/0.20
RoToR	<b>52.8</b>	<b>46.2</b>	<b>53.8</b>	<b>63.5</b>	<b>59.1</b>	<b>65.3</b>	31.8	28.8	35.5	<b>34.2</b>	<b>29.9</b>	<b>37.7</b>	<b>47.4</b>	43.1	49.1
→ stdev. (±)	0.05/0.05/0.02	0.11/0.07/0.08	0.11/0.07/0.08	0.11/0.07/0.08	0.11/0.07/0.08	0.11/0.07/0.08	0.05/0.02/0.09	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.06/0.04/0.07	0.06/0.04/0.07	0.06/0.04/0.07
RoToR-MonoT5	51.6	45.0	52.6	-	-	-	<b>32.4</b>	29.2	36.3	33.0	28.8	36.5	-	-	-
→ stdev. (±)	0.12/0.06/0.10	-	-	-	-	-	0.04/0.02/0.13	0.12/0.09/0.07	0.12/0.09/0.07	0.12/0.09/0.07	0.12/0.09/0.07	0.12/0.09/0.07	-	-	-
RoToR-Freq.	52.5	45.9	53.5	-	-	-	32.3	<b>29.3</b>	36.0	33.8	29.6	37.4	-	-	-
→ stdev. (±)	0.10/0.15/0.11	-	-	-	-	-	0.13/0.16/0.09	0.04/0.00/0.09	0.04/0.00/0.09	0.04/0.00/0.09	0.04/0.00/0.09	0.04/0.00/0.09	-	-	-

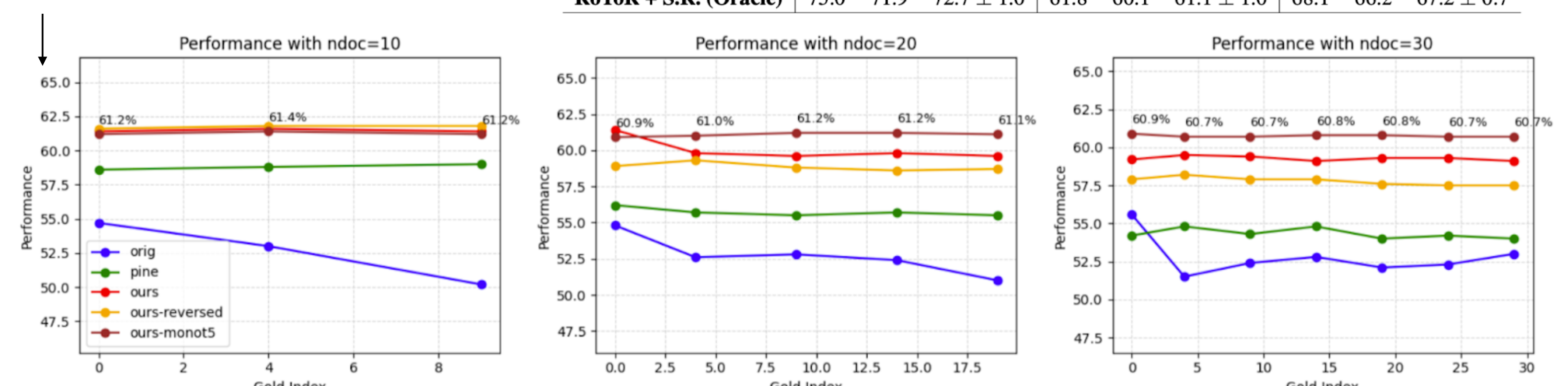
	Llama-3.1-8B-Instr.			Llama-3.1-70B-Instr.			Qwen1.5-4B-Chat			Qwen1.5-7B-Chat			Qwen1.5-72B-Chat		
Method	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1	Acc.	EM	F1
<b>N = 50</b>															
<b>Initial, no shuffling of segments</b>															
Original	50.0	44.0	51.7	62.6	58.5	64.5	31.6	28.6	35.8	31.7	28.0	35.7	42.1	38.7	44.5
PINE	51.6	45.1	52.6	64.1	59.8	65.8	31.6	28.8	35.3	32.0	28.5	35.9	48.0	44.1	49.0
RoToR	52.9	46.0	<b>53.7</b>	<b>64.6</b>	<b>60.0</b>	<b>66.2</b>	<b>32.7</b>	<b>29.6</b>	<b>36.2</b>	<b>34.3</b>	<b>30.1</b>	<b>38.0</b>	<b>48.4</b>	<b>44.3</b>	<b>50.3</b>
RoToR-MonoT5	52.4	45.4	52.8	-	-	-	32.3	29.3	35.9	32.9	28.9	36.6	-	-	-
RoToR-Freq.	<b>53.1</b>	<b>46.4</b>	<b>53.7</b>	-	-	-	32.3	29.2	36.1	33.5	29.5	37.2	-	-	-
<b>After shuffling segments, averaged over 3 seeds</b>															
Original	49.7	43.5	51.0	62.8	58.5	64.5	30.3	27.6	35.0	31.6	27.9	35.5	42.1	38.9	44.7
→ stdev. (±)	0.34/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.37/0.28/0.46	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.26/0.24/0.43	0.30/0.40/0.33	0.30/0.40/0.33	0.30/0.40/0.33
PINE	51.8	45.3	52.7	64.3	59.8	65.9	31.5	28.7	35.3	31.7	28.2	35.7	<b>48.0</b>	<b>44.3</b>	<b>50.0</b>
→ stdev. (±)	0.05/0.07/0.16	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.13/0.04/0.10	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.20/0.18/0.13	0.18/0.20/0.20	0.18/0.20/0.20	0.18/0.20/0.20
RoToR	52.7	45.9	<b>53.5</b>	<b>64.5</b>	<b>60.0</b>	<b>66.1</b>	<b>32.5</b>	<b>29.6</b>	<b>36.1</b>	<b>34.2</b>	<b>30.1</b>	<b>38.0</b>	<b>48.3</b>	<b>44.3</b>	<b>50.3</b>
→ stdev. (±)	0.05/0.09/0.04	0.11/0.06/0.09	0.11/0.06/0.09	0.11/0.06/0.09	0.11/0.06/0.09	0.11/0.06/0.09	0.05/0.02/0.09	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.09/0.07/0.06	0.06/0.04/0.07	0.06/0.04/0.07	0.06/0.04/0.07
RoToR-MonoT5	51.6	45.2	52.8	-	-	-	32.3	29.4	35.9	32.8	28.8	36.5	-	-	-
→ stdev. (±)	0.12/0.06/0.10	-	-	-	-	-	0.16/0.13/0.07	0.16/0.09/0.07	0.16/0.09/0.07	0.16/0.09/0.07	0.16/0.09/0.07	0.16/0.09/0.07	-	-	-
RoToR-Freq.	<b>53.1</b>	<b>46.4</b>	<b>53.7</b>	-	-	-	32.4	29.3	<b>36.1</b>	33.7	29.6	37.4	-	-	-
→ stdev. (±)	0.02/0.07/0.03	-	-	-	-	-	0.09/0.04/0.06	0.04/0.16/0.22	0.04/0.16/0.22	0.04/0.16/0.22	0.04/0.16/0.22	0.04/0.16/0.22	-	-	-

- Top-30 and Top-50 knowledge triples per query
- Test before / after shuffling segments to see robustness
- RoToR obtains lower stdev (better stability) + higher performance than PINE
- Trend persists for > 70B model variants

- Single use of Order-invariant models fail, but selective routing restores accuracy & improves performance.

- On Llama-3.1-8B-Inst.
- Original model fluctuates performance, while ours (RoToR) maintains stable

## LitM



## MMLU

	Llama-3.1-8B-Instruct			Qwen1.5-4B-Chat			Qwen1.5-7B-Chat		
Method	Init.	Rev.	Avg.	Init.	Rev.	Avg.	Init.	Rev.	Avg.
Orig.	68.3	64.8	65.5 ± 1.0	53.6	51.9	52.6 ± 0.6	<b>60.1</b>	56.6	58.6 ± 0.9
PCW	57.0	55.1	56.1 ± 1.1	-	-	-	-	-	-
Set-Based Prompting	31.1	33.0	31.6 ± 0.8	-	-	-	-	-	-
PINE	64.8	63.3	63.6 ± 0.7	50.5	49.3	49.4 ± 0.5	57.0	54.4	55.8 ± 0.9
RoToR	63.2	62.6	62.8 ± 0.7	49.6	47.7	48.3 ± 0.7	56.5	55.8	56.2 ± 0.6
→ + S.R.	<b>68.5</b>	65.1	65.7 ± 0.9	53.7	51.8	<b>52.6 ± 0.6</b>	<b>60.1</b>	<b>57.4</b>	<b>58.8</b>