

# Persistent Topological Features in Large Language Models

**Yuri Gardinazzi<sup>1,2</sup>, Karthik Viswanathan<sup>1,3</sup>, Giada Panerai<sup>1,2</sup>**  
Alessio Ansuini<sup>1</sup>, Alberto Cazzaniga<sup>1</sup>, Matteo Biagetti<sup>†,1</sup>

<sup>1</sup>AREA Science Park, <sup>2</sup>University of Trieste, <sup>3</sup>University of Amsterdam



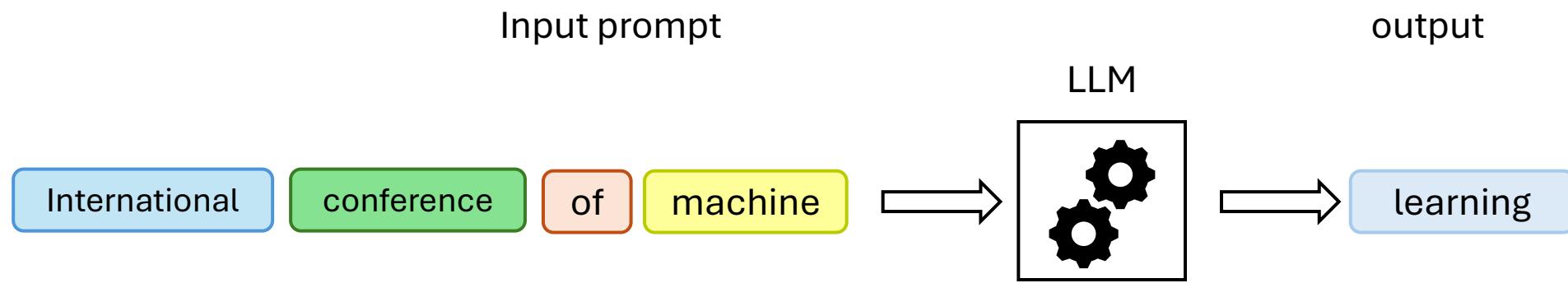
ICML 2025

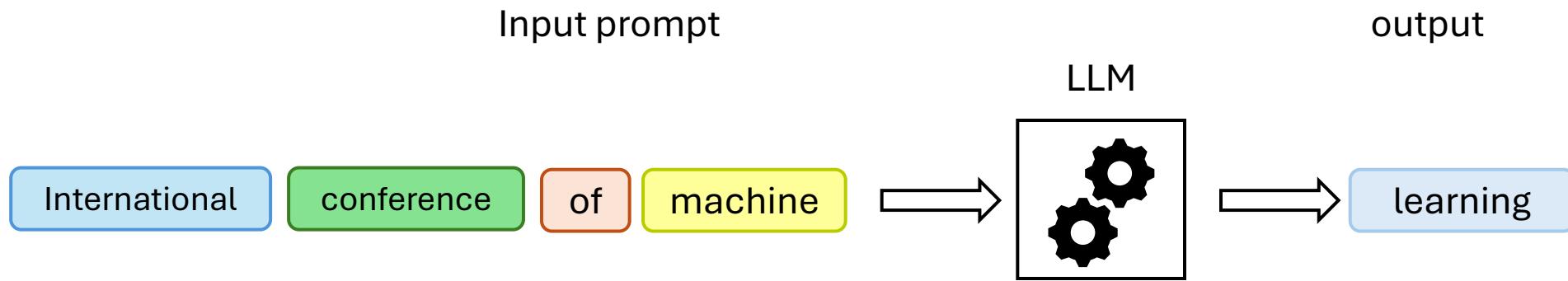


UNIVERSITÀ  
DEGLI STUDI  
DI TRIESTE

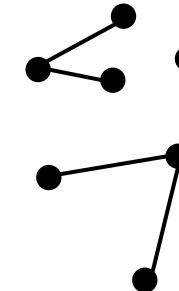


UNIVERSITY OF AMSTERDAM

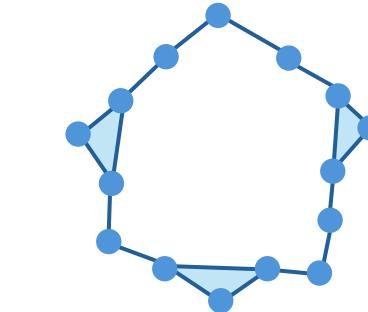




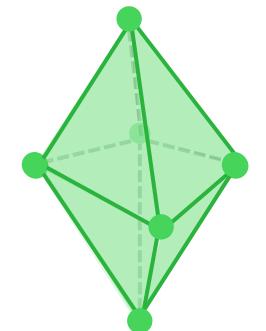
## Topological Features



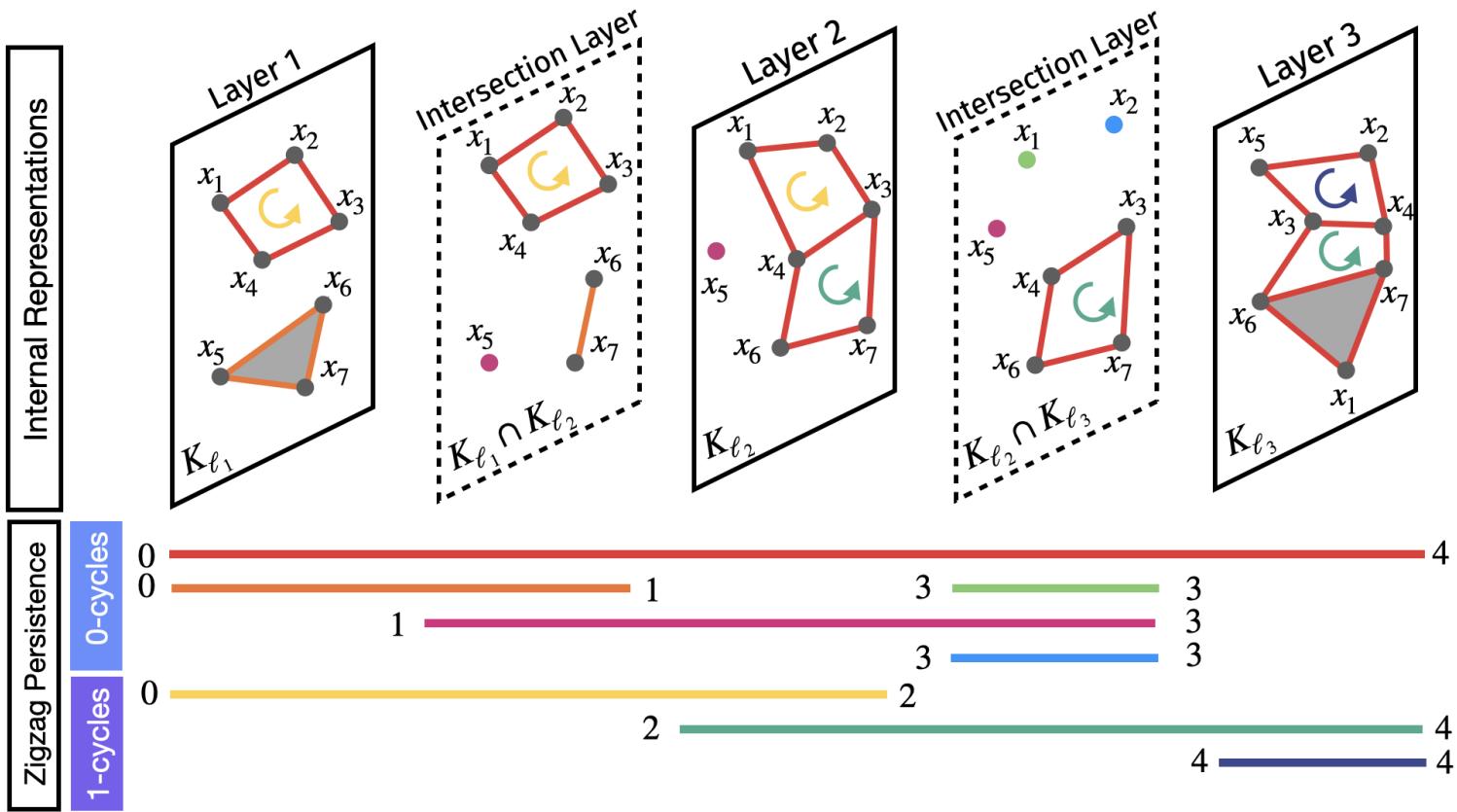
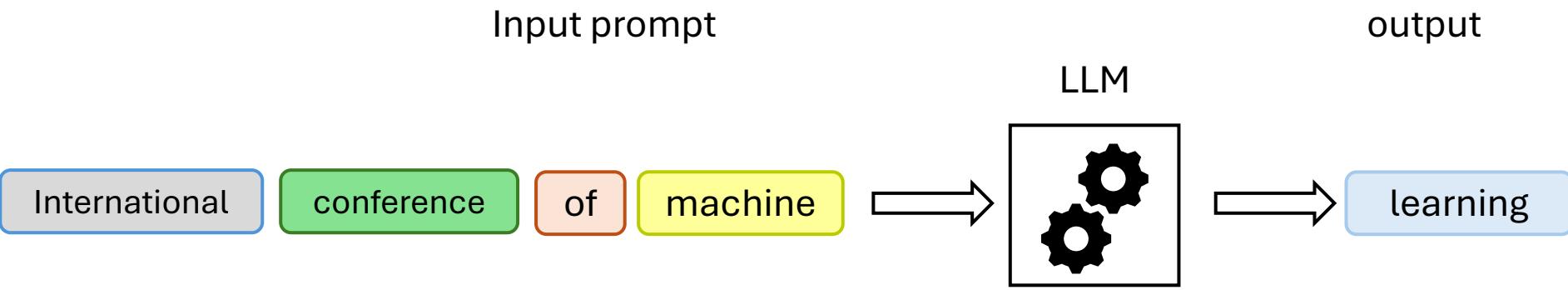
Connected components



Loops



Voids



# Input prompt

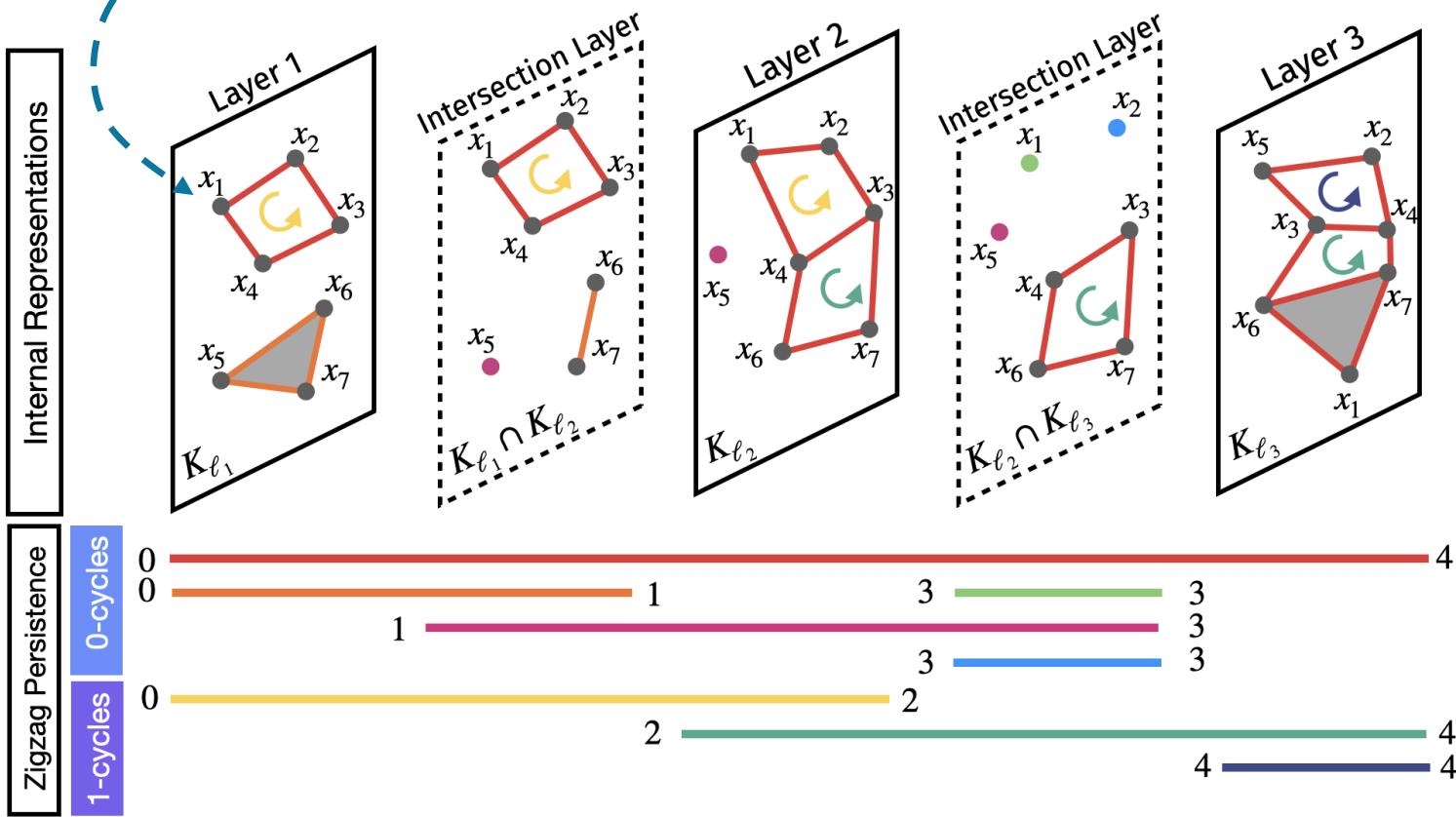
International

conference

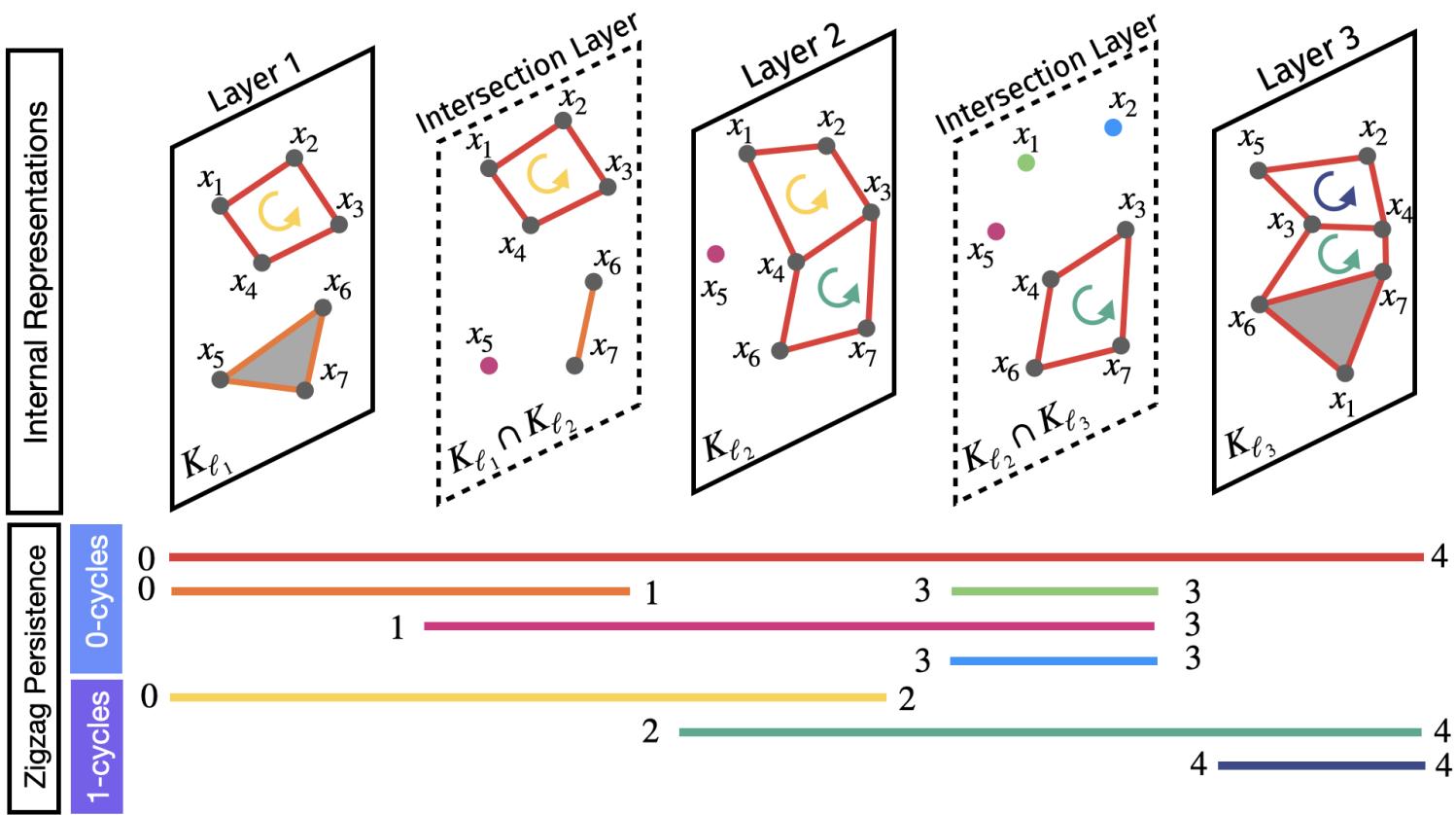
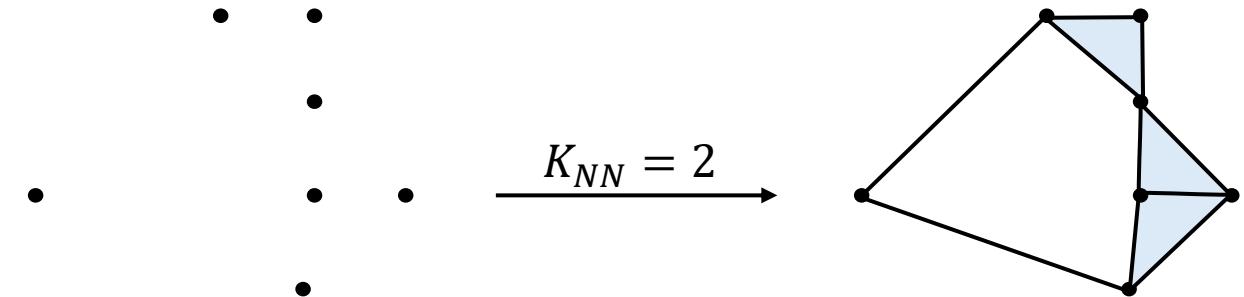
of

machine

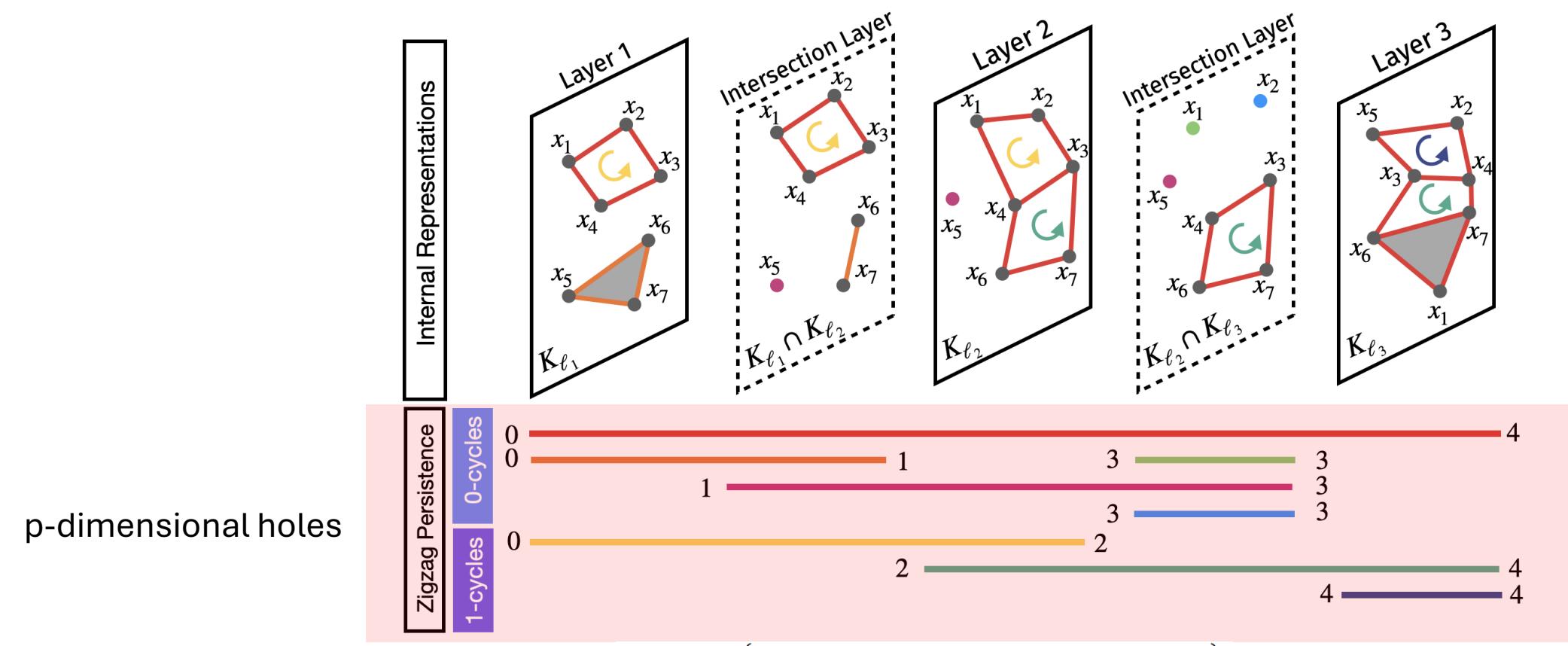
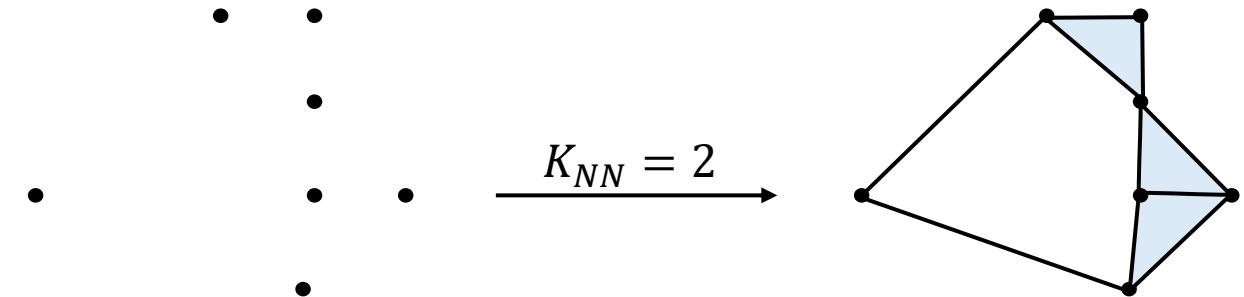
$$\in \mathbb{R}^d$$



# Simplicial Complex



# Simplicial Complex



# Topological Descriptors

Births' relative frequency

$$B_p(\ell) = \frac{\sum_{\ell_i} \omega(\ell, \ell_i) \widehat{PI}_p(\ell, \ell_i)}{\sum_{\ell_i} \omega(\ell, \ell_i) \sum_{\ell_i} \widehat{PI}_p(\ell, \ell_i)}$$

$$\omega(\ell, \ell_i) = |\ell - \ell_i|^\alpha$$

$$\mathcal{Z}_p(\ell_1, \ell_2) = \frac{\sum_{\ell_1 \leq M_1, \ell_2 > M_2} \widehat{PI}_p(\ell_1, \ell_2)}{\beta_p(\ell_1)}$$

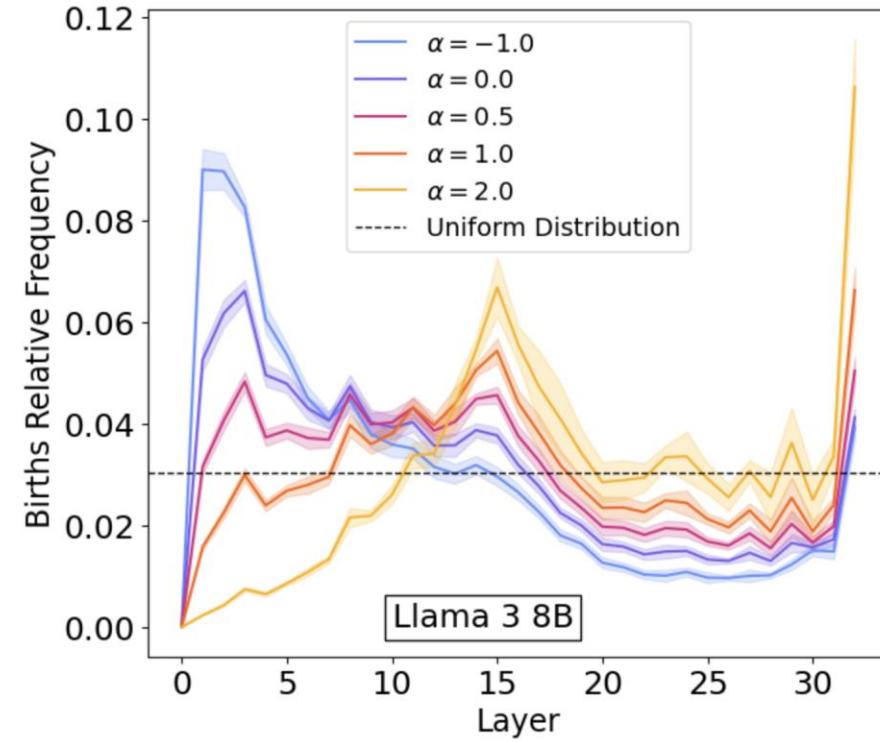
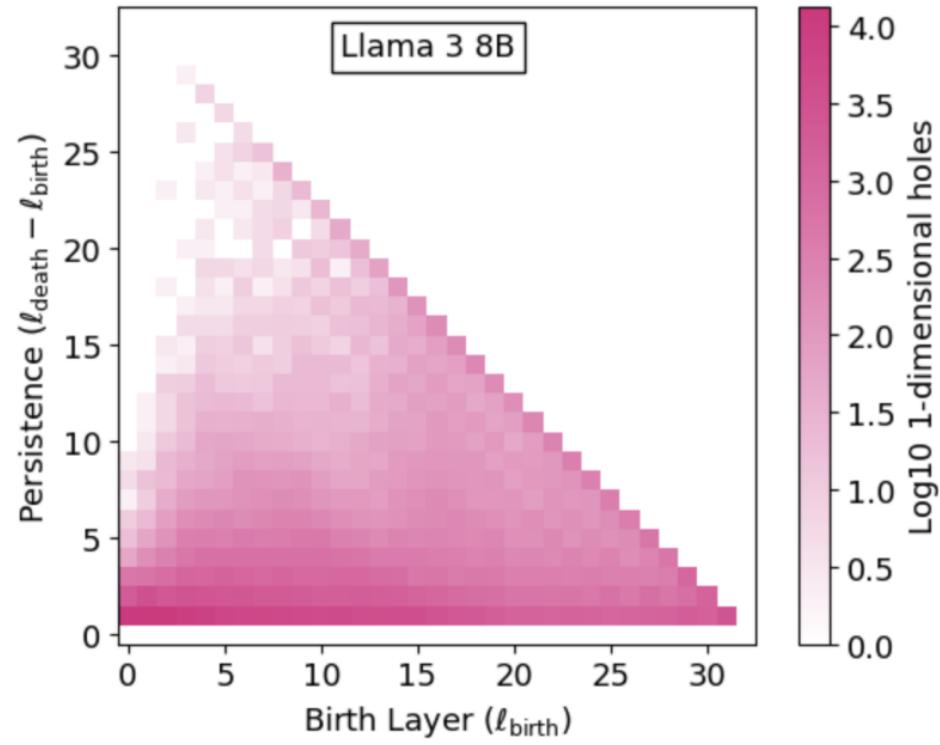
Weighed Inter-Layer Persistence

$$\bar{\mathcal{Z}}_p(\ell) = \frac{\sum_{\ell_i=1}^{N_{\text{layers}}} \omega(\ell, \ell_i) \mathcal{Z}_p(\ell, \ell_i)}{\sum_{\ell_i=1}^{N_{\text{layers}}} \omega(\ell, \ell_i)}$$

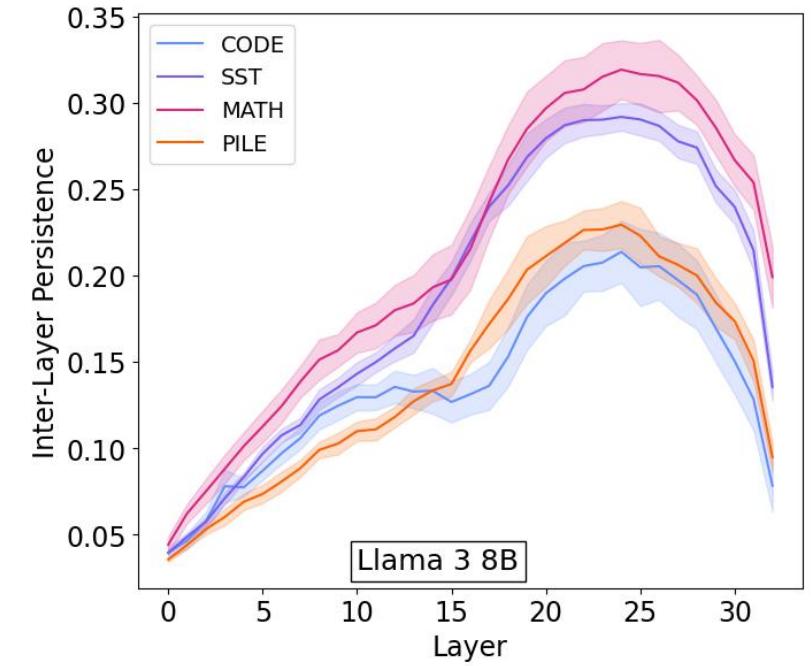
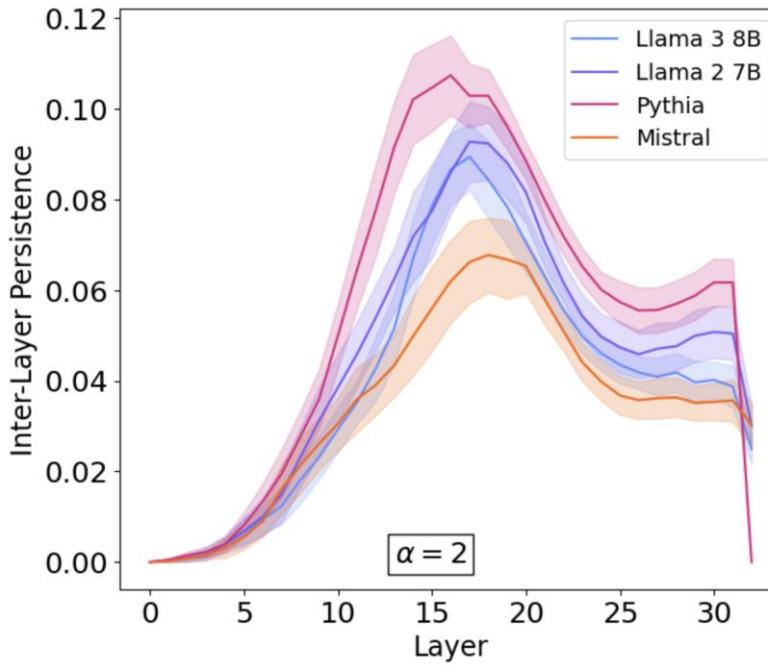
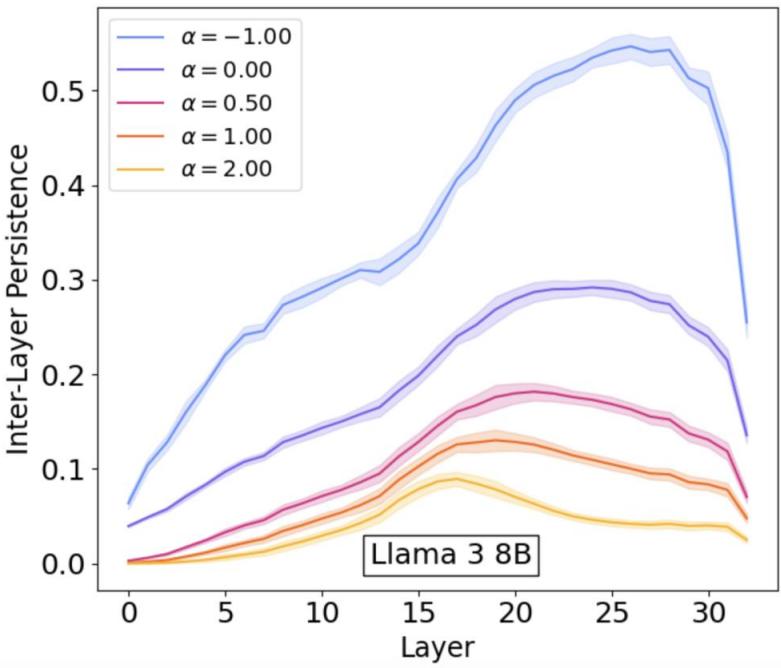
$$M_1 = \min(\ell_1, \ell_2)$$

$$M_2 = \max(\ell_1, \ell_2)$$

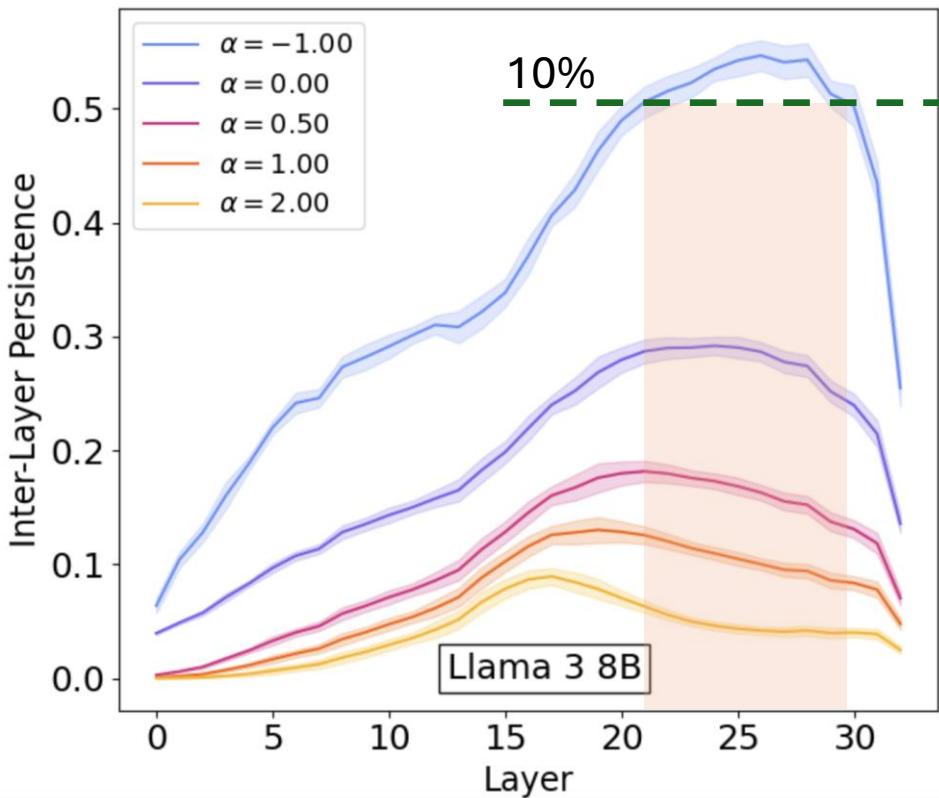
# Results



# Results



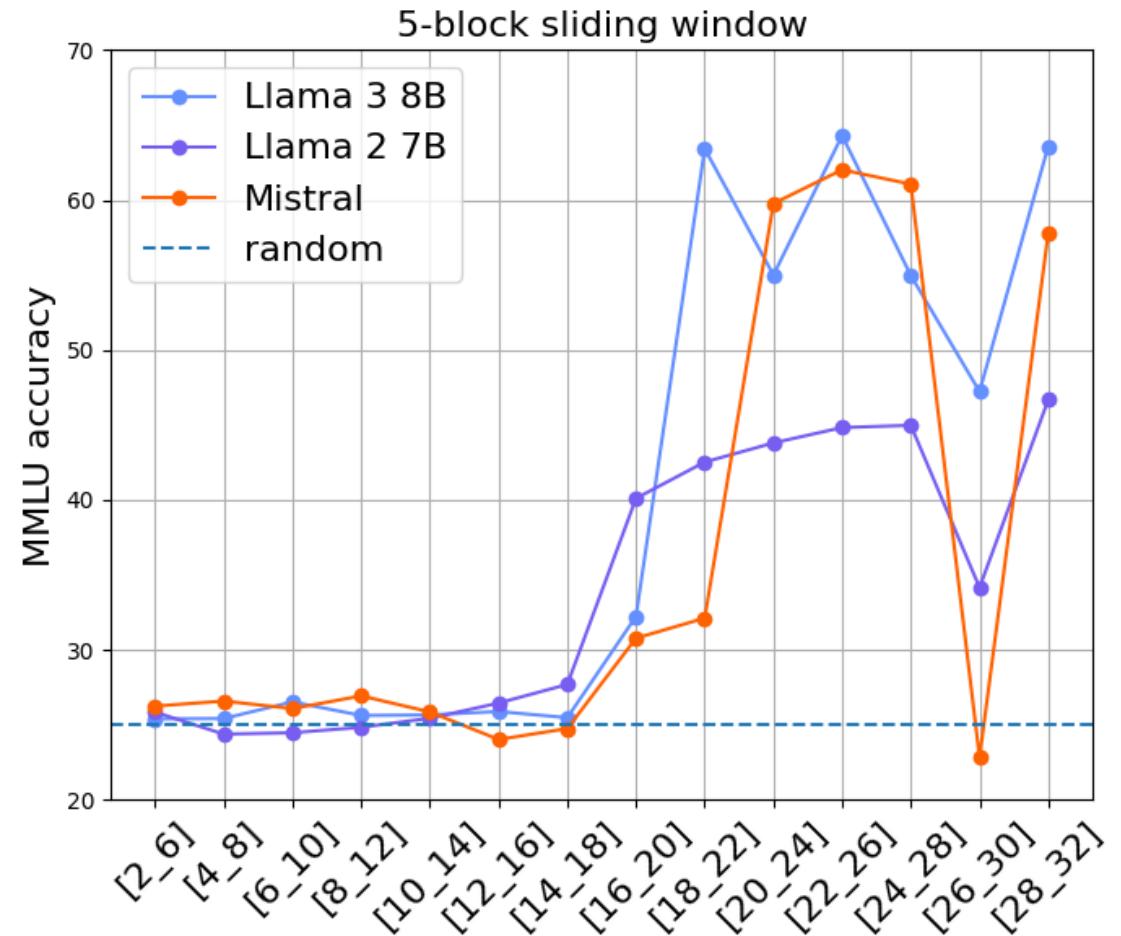
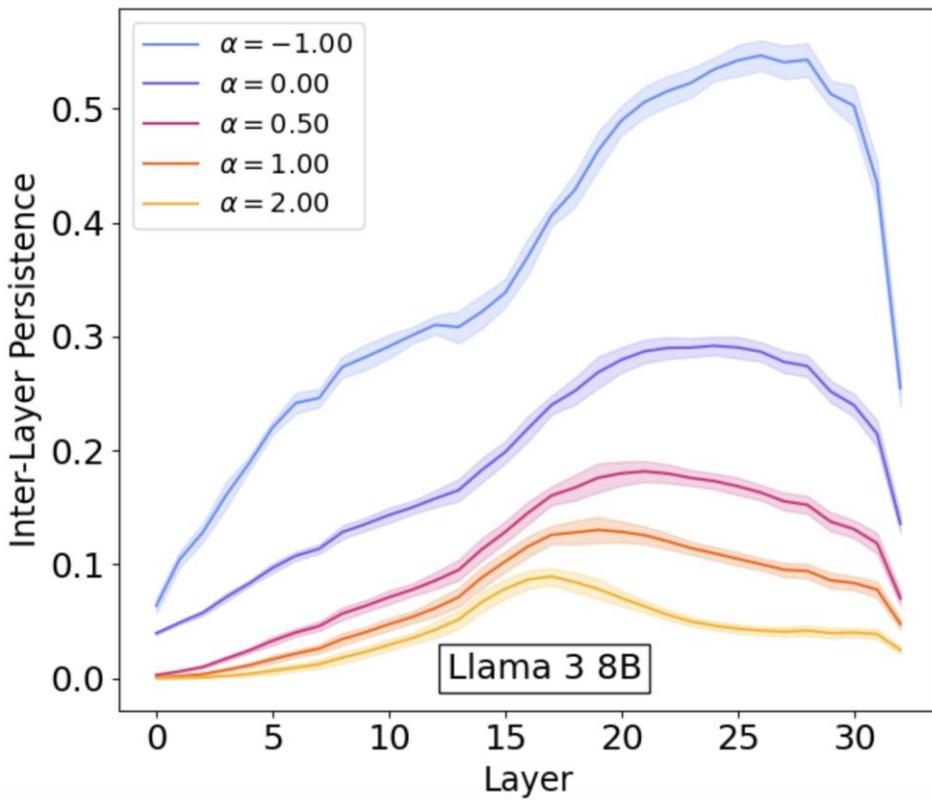
# Results



Models	MMLU			HellaSwag			WinoGrande		
	Full	This work	Other works	Full	This work	Other works	Full	This work	Other works
Llama 2	45.74	37.38	<b>43.95</b>	58.54	<b>44.71</b>	42.78	74.43	<b>68.67</b>	67.72
Llama 3	65.07	<b>53.44</b>	<b>53.44</b>	61.37	<b>41.60</b>	<b>41.60</b>	77.10	<b>70.00</b>	<b>70.00</b>
Mistral 7B	62.40	<b>53.17</b>	38.20	62.83	<b>36.67</b>	34.45	77.35	<b>66.50</b>	63.76
Pythia	-	-	-	49.70	31.43	<b>34.96</b>	63.30	55.71	<b>58.09</b>

Other works: **Gromov et al. (2024)**, **Men et al. (2024)**

# Results



# Conclusions

- We apply TDA to interpret LLMs behaviors.
- With ZigZag we study **trajectories** of point clouds that evolve in time (or through layers).
- We can distinguish different **phases** of prompt processing.
- We can measure the **rearrangement** of points through different layers.

# Find us at ICML 2025!

- Poster Session: Wed 16 Jul 4:30 p.m. PDT — 7 p.m. PDT
- Email: [yuri.gardinazzi@areasciencepark.it](mailto:yuri.gardinazzi@areasciencepark.it)