

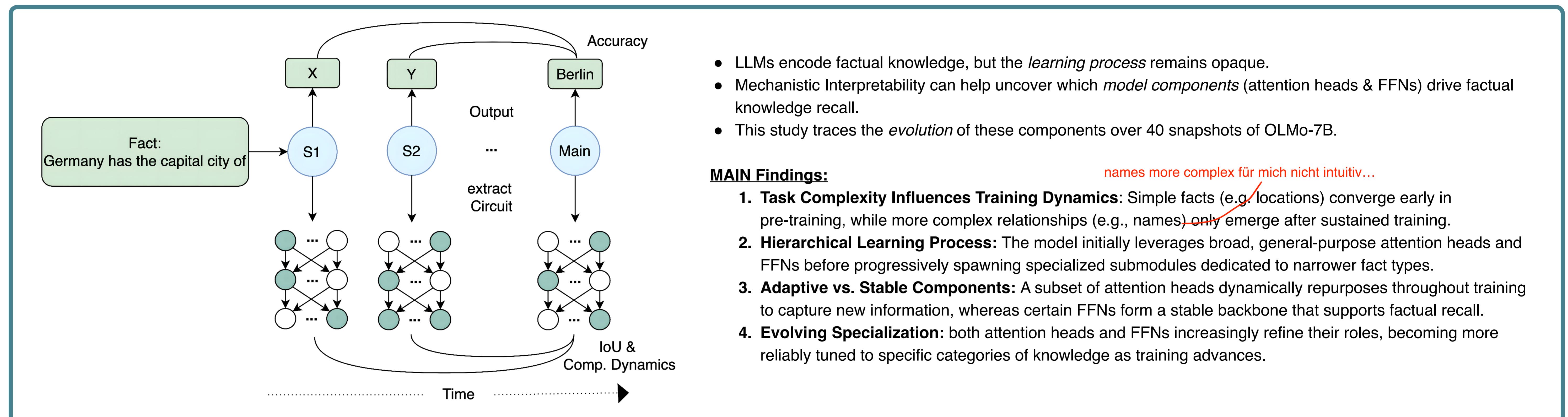
Time Course MechInterp: Analyzing the Evolution of Components and Knowledge in Large Language Models

ACL 2025
VIENNA

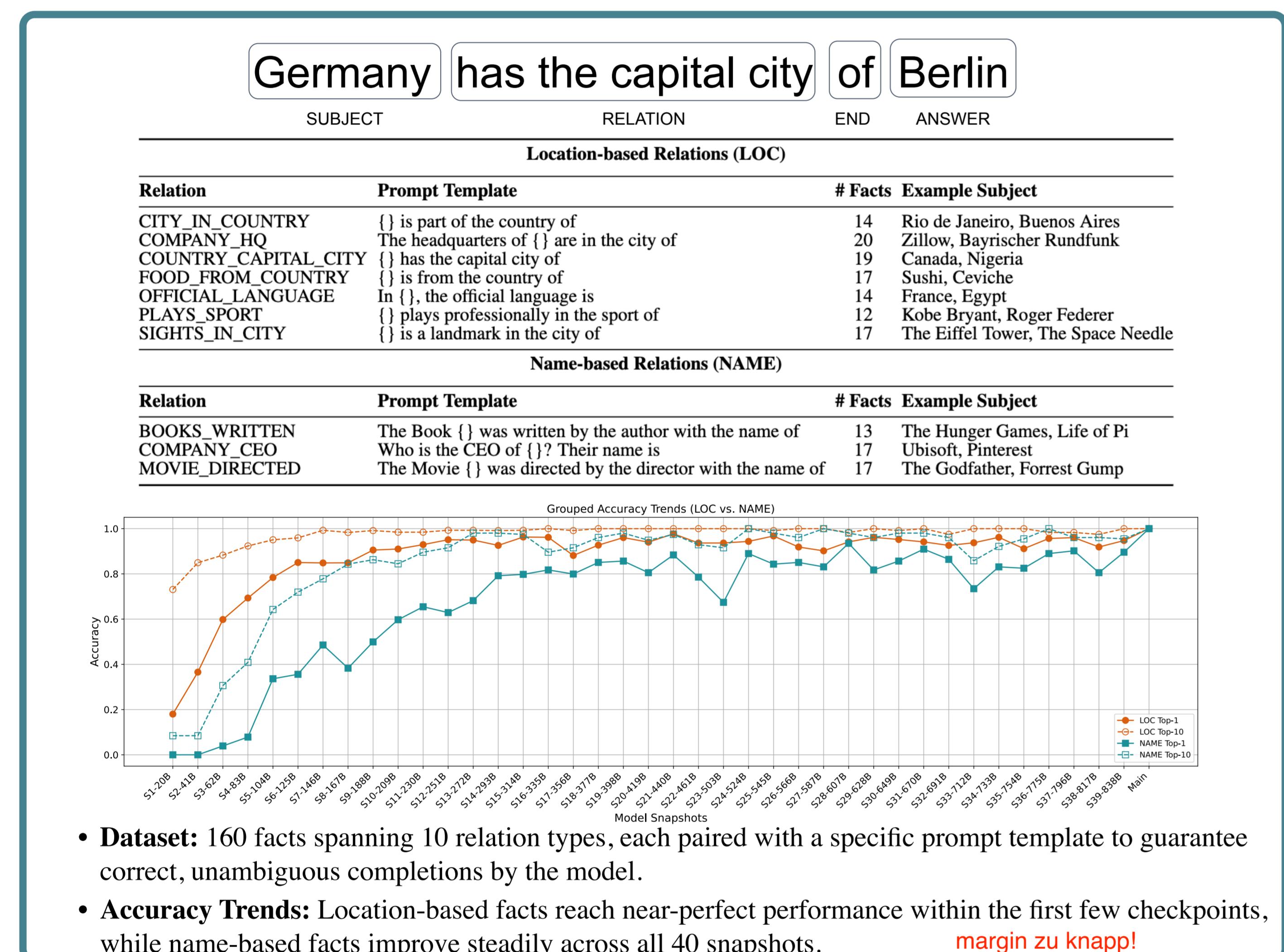
mcmL
Munich Center for Machine Learning

Ahmad Dawar Hakimi, Ali Modarressi, Philipp Wicke, Hinrich Schütze

How does factual knowledge emerge during LLM pretraining?



Dataset Construction



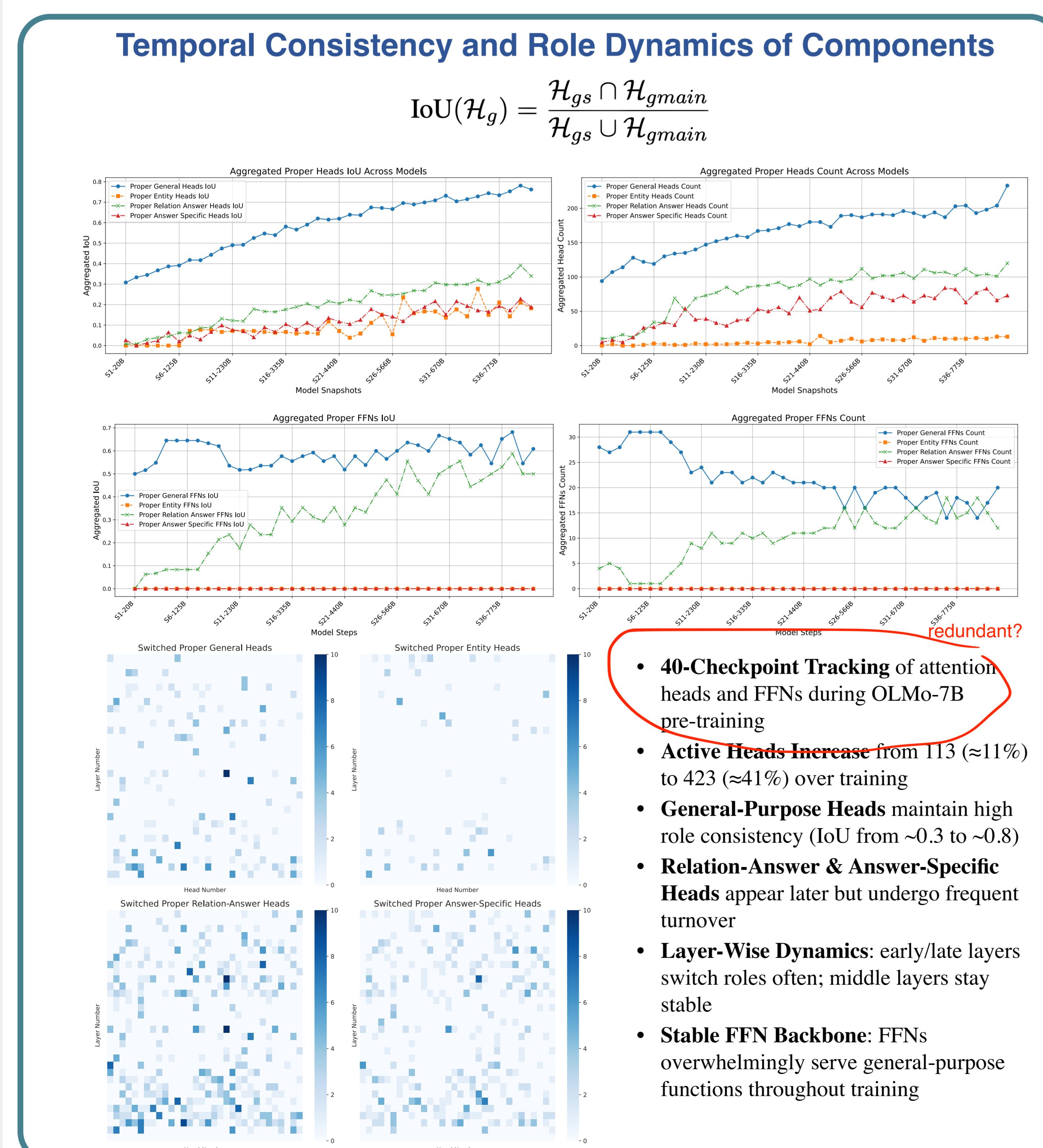
Model Component Roles

Role Score:		Hierarchical Proper Role:
General:	$c_s^g = \frac{\sum_{r \in R} \sum_{f \in r} c_{srf}(T_g)}{\sum_{r \in R} \sum_{f \in r} 1}$	$\mathcal{H}_g = \mathcal{J}_g$
Entity:	$c_s^e = \frac{\sum_{r \in R} \sum_{f \in r} c_{srf}(T_e)}{\sum_{r \in R} \sum_{f \in r} 1}$	$\mathcal{H}_e = \mathcal{J}_e - \mathcal{J}_g$
Relation-Answer:	$c_s^r = \frac{\sum_{f \in r} c_{srf}(T_a)}{\sum_{f \in r} 1}$	$\mathcal{H}_r = \mathcal{J}_r - \mathcal{J}_e - \mathcal{J}_g$
Fact-Answer Specific:	$c_s^f = c_{srf}^f(T_a)$	$\mathcal{H}_f = \mathcal{J}_f - \mathcal{J}_r - \mathcal{J}_e - \mathcal{J}_g$

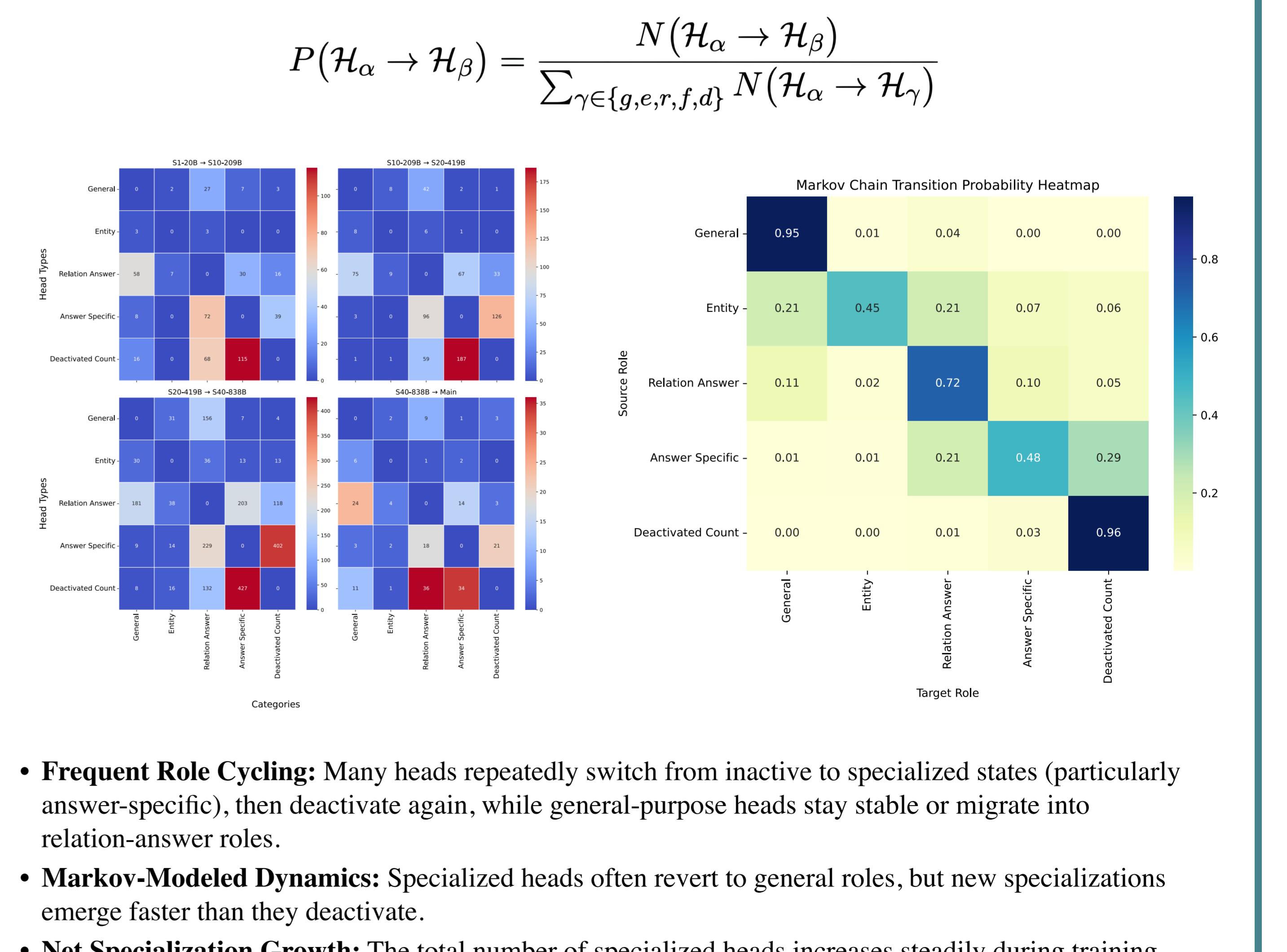
Deactivated Components: $\mathcal{H}_d = \mathcal{C}(\mathcal{H}_g \cup \mathcal{H}_e \cup \mathcal{H}_r \cup \mathcal{H}_f)$

- Extract per-subtoken circuits using Information Flow Routes (Ferrando & Voita, 2024)
- We compute activation scores $c_s^g, c_s^e, c_s^r, c_s^f$ for each component at each snapshot, using subtoken sets T_g, T_e, T_a , and a threshold $\theta = 0.1$.
- Then by successive differencing of cumulative importance sets $J_g \rightarrow J_e \rightarrow J_r \rightarrow J_f$, we obtain non-overlapping proper role sets H_g, H_e, H_r and H_f , with H_d capturing all remaining deactivated components.

How do Components Evolve?



Dynamic Specialization and Generalization of Attention Heads



Code



Paper



Poster

