

# COLLAGE: Collaborative Human-Agent Interaction Generation using Hierarchical Latent Diffusion and Language Models

Divyanshu Daiya<sup>1</sup>, Damon Conover<sup>2</sup>, Aniket Bera<sup>1</sup>

<sup>1</sup> Purdue University, <sup>2</sup> DEVCOM Army Research Laboratory

## Motivation

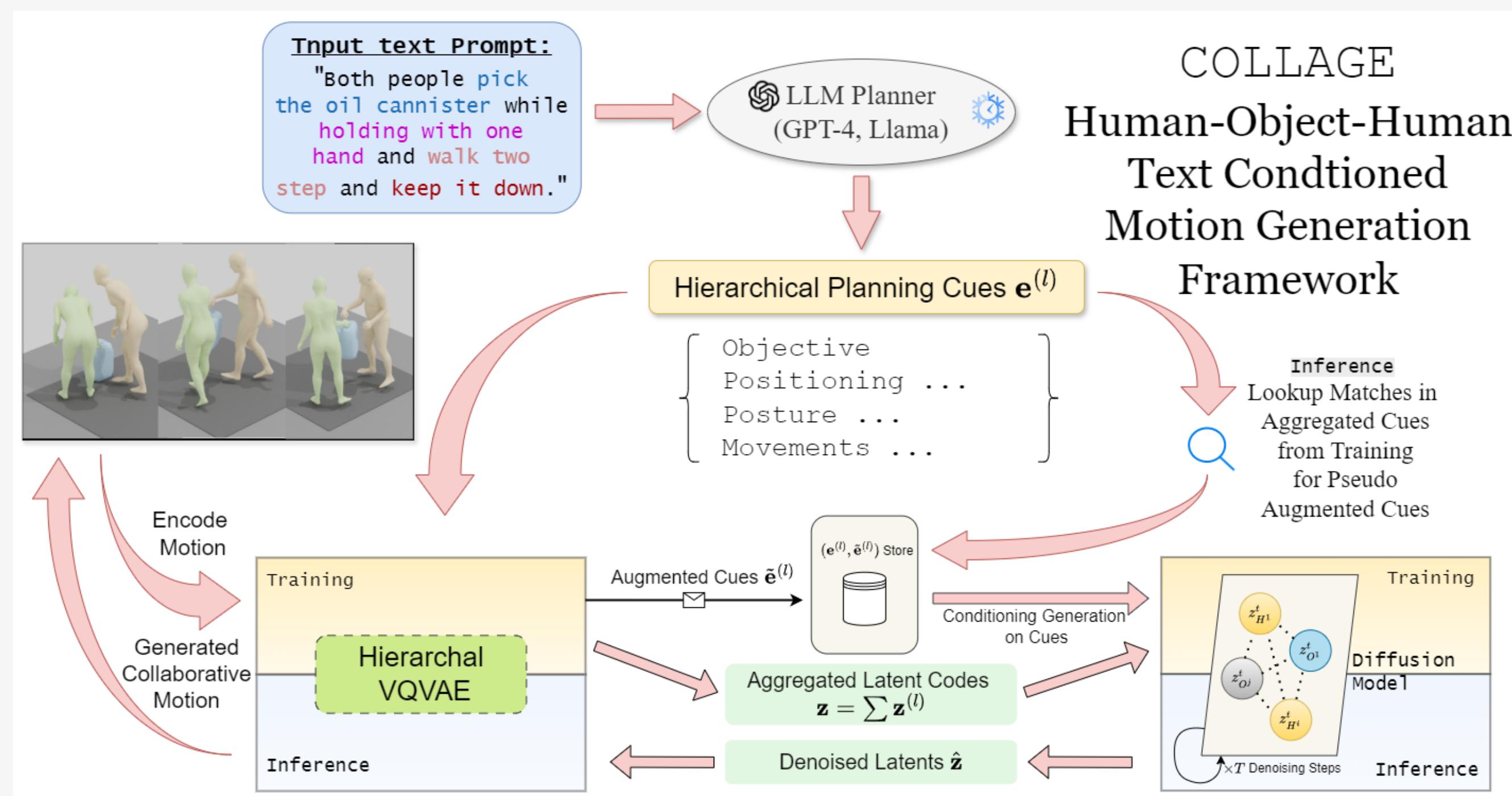
**Challenges.** Generating collaborative (Human–Object–Human) motions is difficult because of:

- Complex multi-agent planning,
- Limited datasets that capture fine-grained synergy,
- Need for fine control over interactions.

### Key Contributions.

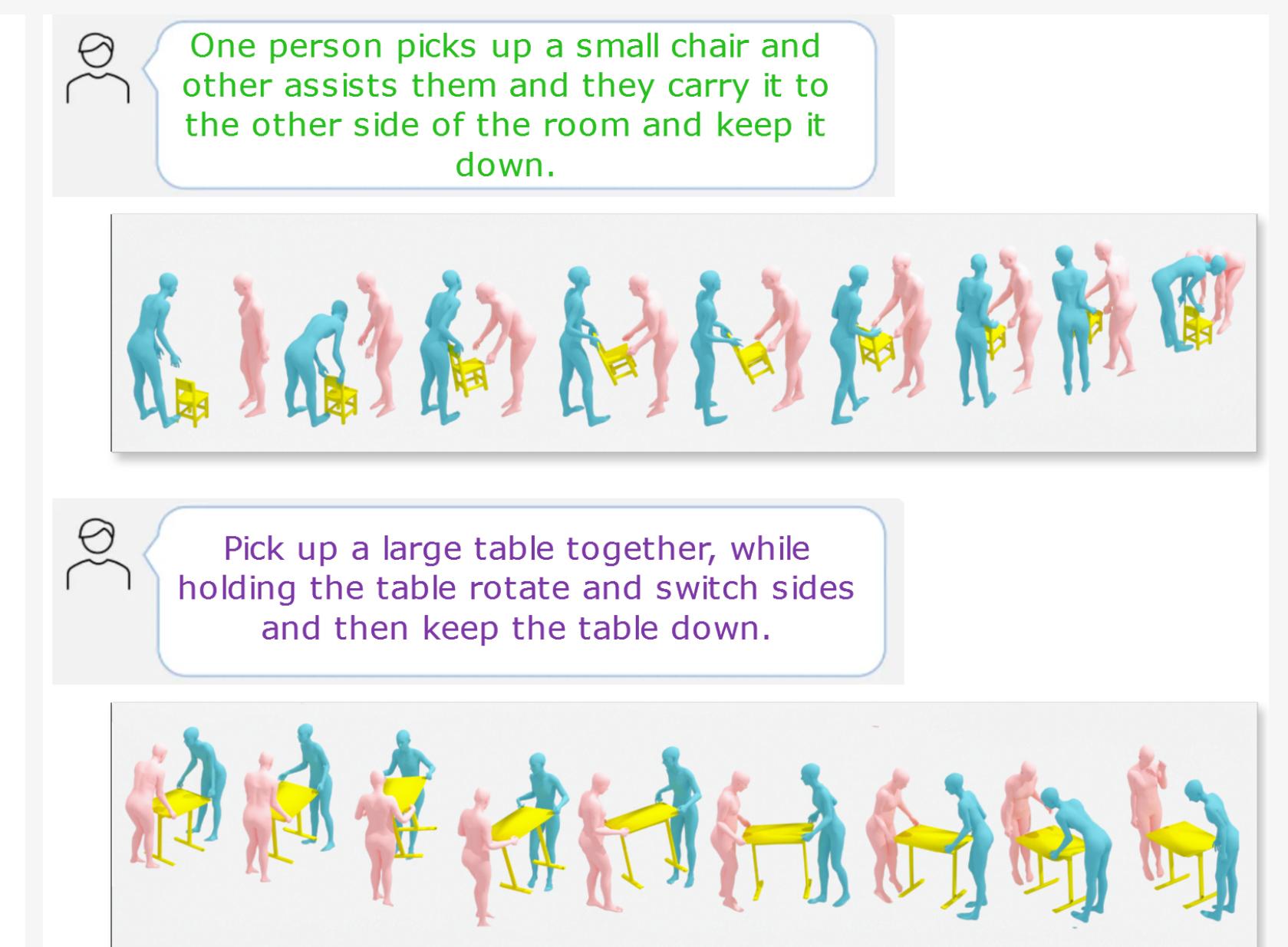
- **LLM-guided hierarchical VQ-VAE** supplies a structured, coarse-to-fine action script, directly taming the complexity of multi-agent planning.
- **Text ↔ code semantic anchors** learned via contrastive training enrich sparse data with strong language priors, compensating for limited fine-grained synergy examples.
- **Anchor-conditioned diffusion** refines motions in codebook space within **5–55 steps**, delivering precise, high-fidelity control and achieving SOTA generative performance on CORE4D and InterHuman datasets.

**Goal:** Generate realistic multi-actor + object collaboration from text or geometry.



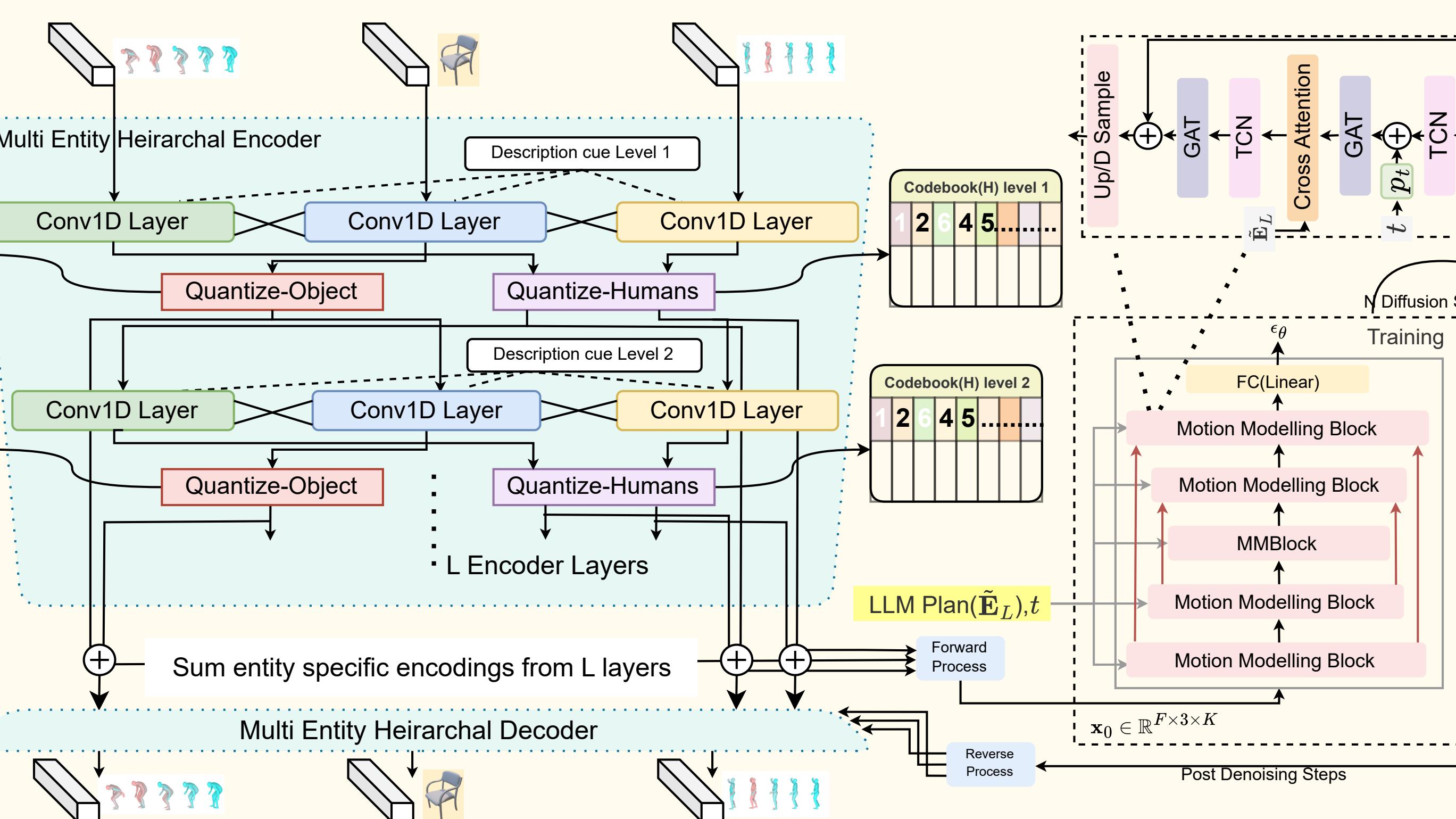
**High-Level Insight:** LLM-based semantics + hierarchical latents yield coherent, controllable multi-agent interactions.

## Overview



## Methodology

**Hierarchical Planning Cues from LLM for Collaborative Task**  
This structure guides two agents working together to carry a chair across the floor. Instructions span multiple levels, from high-level goals to fine-grained control.  
Layer 1: Overall Objective  
Define the key goal at the highest level of the motion task:  
Goal: Both Person1 and Person2 collaborate to carry a chair from the left side of the room to the right.  
Layer 2: Initial Positioning and Approach  
Positioning of the agents and their initial movement towards the object:  
Person1 approaches from the left side of the chair.  
Person2 approaches from the right side.  
Both face the chair directly, establishing initial contact.  
Layer 3: Contact Points and Gripping Mechanics  
Define how each person grips the object:  
Person1 grasps the back of the chair with their left hand.  
Person2 grasps the seat of the chair firmly with their right hand.  
The chair is securely held by both individuals.  
Layer 4: Upper Body Posture and Coordination  
Details upper body posture for balance:  
Person1 maintains a level center of gravity to evenly distribute the weight.  
Person2 leans slightly forward to maintain control.  
Both maintain upright posture for proper balance.  
Layer 5: Lower Body Movements and Step Synchronization  
Instructions for synchronized movement:  
Person1 takes careful, deliberate steps forward.  
Person2 synchronizes steps with Person1 to ensure balance.  
Both shift weight slightly forward with each step for smooth motion.  
Layer 6: Fine-Grained Adjustments and Limb Coordination  
Small adjustments to avoid obstacles and maintain coordination:  
Person1 extends left arm to guide the chair direction.  
Person2 flexes right arm to adjust the chair's angle as needed.  
Both make minor adjustments to avoid obstacles and maintain smooth path.



### Latent Generation & Planning.

- A Graph U-Net built from TCN + GAT blocks models multi-actor synergy.
- Training minimises the **noise-prediction loss**  $\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{E}_L)\|^2$ .
- **Time-modulated weights** emphasise broad constraints early and fine detail late.
- Converges in just **5–55 DDIM steps**.

### Incorporating Language.

- **Contrastive links** tie textual prompts to the learned codebook, guaranteeing semantic relevance.
- **Cross-attention** fuses text and motion signals inside every TCN/GAT layer.
- **Hierarchical GPT cues** such as “Approach → Grasp → Lift” guide generation step-by-step along the diffusion timeline.

### Hierarchical VQ-VAE.

- Each layer encodes a different motion scale, and **LLM text cues** are injected at every level to steer semantics in real time.
- **Cross-code / text alignment** binds discrete latent codes to sentence embeddings, giving robust semantic anchors.
- Those anchors guide **anchor-conditioned diffusion** at test-time, so motions stay text-consistent while requiring markedly fewer denoise steps.
- A **compact loss mix**—reconstruction, contact, penetration and alignment—jointly enforces physical plausibility and semantic fidelity.
- The resulting latents are **disentangled and stable**, enabling smooth multi-agent synergy.

### Iterative Refinement.

- **LLM sub-actions** are injected at each partial denoise step, acting as intermediate waypoints.
- Codes **progressively align** to those sub-actions while Gaussian noise shrinks, tightening the motion trajectory.
- The procedure preserves consistent agent – agent - object collaboration even on long, complex tasks.

## Experimental Results

**Datasets:** CORE-4D & InterHuman; **Metrics:** R-Precision (Top-1/2/3), FID, Diversity, Multimodality, MM-Distance, Joint Position Error (RR.J<sub>e</sub>), Vertex Position Error (RR.V<sub>e</sub>), Contact Accuracy (C<sub>acc</sub>).

### Text-Conditioned Generation on CORE4D and InterHuman

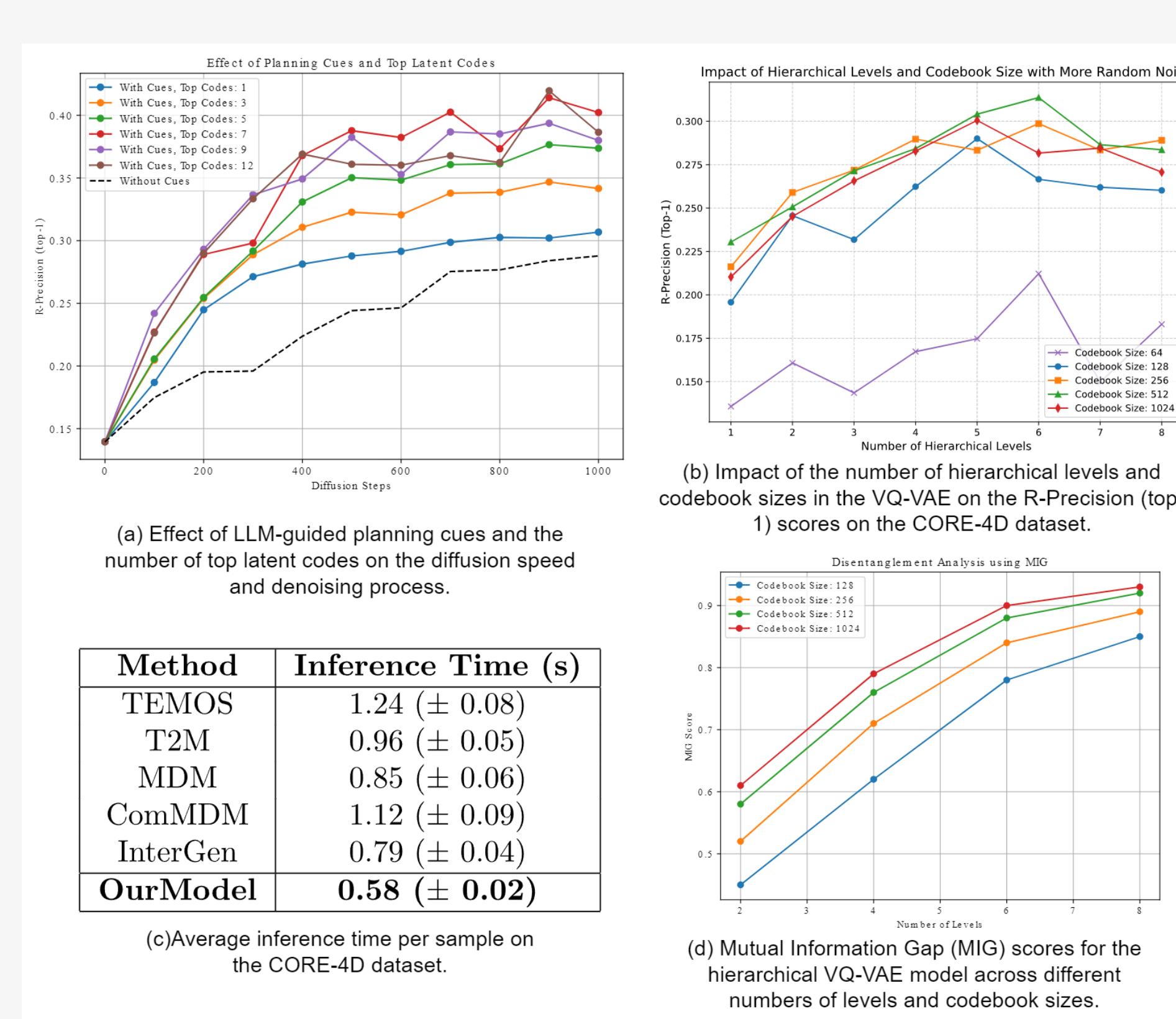
Methods(CORE-4D)	R Precision↑			FID↓	MM Dist.↓	Diversity→	MModality↑	Methods(InterHuman)	R Precision↑			FID↓	MM Dist.↓	Diversity→	MModality↑
	Top 1	Top 2	Top 3						Top 1	Top 2	Top 3				
Real	0.312 <sup>±0.007</sup>	0.587 <sup>±0.006</sup>	0.673 <sup>±0.006</sup>	0.005 <sup>±0.0005</sup>	4.124 <sup>±0.019</sup>	8.151 <sup>±0.091</sup>	-	Real	0.459 <sup>±0.008</sup>	0.610 <sup>±0.009</sup>	0.701 <sup>±0.008</sup>	0.273 <sup>±0.007</sup>	3.755 <sup>±0.008</sup>	7.948 <sup>±0.064</sup>	-
TEMOS	0.065 <sup>±0.006</sup>	0.179 <sup>±0.006</sup>	0.211 <sup>±0.005</sup>	9.214 <sup>±0.0758</sup>	8.536 <sup>±0.019</sup>	4.671 <sup>±0.091</sup>	0.510 <sup>±0.052</sup>	TEMOS	0.224 <sup>±0.010</sup>	0.316 <sup>±0.013</sup>	0.450 <sup>±0.018</sup>	17.375 <sup>±0.043</sup>	6.342 <sup>±0.015</sup>	6.939 <sup>±0.071</sup>	0.533 <sup>±0.014</sup>
T2M	0.195 <sup>±0.003</sup>	0.141 <sup>±0.005</sup>	0.267 <sup>±0.002</sup>	11.258 <sup>±0.064</sup>	5.867 <sup>±0.013</sup>	2.738 <sup>±0.076</sup>	1.672 <sup>±0.041</sup>	T2M	0.238 <sup>±0.012</sup>	0.325 <sup>±0.012</sup>	0.464 <sup>±0.014</sup>	13.769 <sup>±0.072</sup>	5.731 <sup>±0.018</sup>	7.046 <sup>±0.022</sup>	1.387 <sup>±0.076</sup>
MDM	0.163 <sup>±0.013</sup>	0.257 <sup>±0.010</sup>	0.348 <sup>±0.008</sup>	9.671 <sup>±0.0629</sup>	10.219 <sup>±0.020</sup>	7.395 <sup>±0.090</sup>	3.526 <sup>±0.074</sup>	MDM(GRU)	0.168 <sup>±0.009</sup>	0.279 <sup>±0.009</sup>	0.361 <sup>±0.010</sup>	10.228 <sup>±0.025</sup>	6.951 <sup>±0.151</sup>	3.170 <sup>±0.046</sup>	2.356 <sup>±0.080</sup>
ComMDM	0.187 <sup>±0.005</sup>	0.256 <sup>±0.007</sup>	0.301 <sup>±0.007</sup>	9.217 <sup>±0.0727</sup>	7.541 <sup>±0.023</sup>	5.367 <sup>±0.080</sup>	0.721 <sup>±0.065</sup>	ComMDM	0.229 <sup>±0.010</sup>	0.334 <sup>±0.008</sup>	0.466 <sup>±0.012</sup>	7.069 <sup>±0.054</sup>	6.219 <sup>±0.021</sup>	7.244 <sup>±0.038</sup>	1.829 <sup>±0.052</sup>
InterGen	0.206 <sup>±0.007</sup>	0.312 <sup>±0.008</sup>	0.401 <sup>±0.008</sup>	7.217 <sup>±0.2321</sup>	10.251 <sup>±0.017</sup>	6.162 <sup>±0.223</sup>	3.402 <sup>±0.063</sup>	InterGen	0.371 <sup>±0.010</sup>	0.515 <sup>±0.012</sup>	0.624 <sup>±0.010</sup>	5.918 <sup>±0.079</sup>	5.108 <sup>±0.014</sup>	7.387 <sup>±0.029</sup>	2.141 <sup>±0.063</sup>
COLLAGUE	0.229 <sup>±0.008</sup>	0.332 <sup>±0.009</sup>	0.435 <sup>±0.009</sup>	6.890 <sup>±0.219</sup>	5.526 <sup>±0.016</sup>	7.373 <sup>±0.237</sup>	3.589 <sup>±0.066</sup>	COLLAGUE	0.383 <sup>±0.005</sup>	0.547 <sup>±0.009</sup>	0.657 <sup>±0.006</sup>	4.987 <sup>±0.266</sup>	4.993 <sup>±0.012</sup>	7.515 <sup>±0.214</sup>	2.872 <sup>±0.057</sup>
w/o Hierarchy	0.201 <sup>±0.007</sup>	0.309 <sup>±0.008</sup>	0.411 <sup>±0.008</sup>	7.452 <sup>±0.2381</sup>	5.588 <sup>±0.018</sup>	6.995 <sup>±0.224</sup>	3.209 <sup>±0.058</sup>	w/o Hierarchy	0.353 <sup>±0.009</sup>	0.521 <sup>±0.010</sup>	0.632 <sup>±0.009</sup>	5.543 <sup>±0.2154</sup>	5.048 <sup>±0.015</sup>	7.137 <sup>±0.201</sup>	2.497 <sup>±0.061</sup>
w/o LLM	0.208 <sup>±0.007</sup>	0.315 <sup>±0.008</sup>	0.419 <sup>±0.008</sup>	7.235 <sup>±0.2305</sup>	5.561 <sup>±0.017</sup>	6.549 <sup>±0.230</sup>	3.152 <sup>±0.063</sup>	w/o LLM	0.362 <sup>±0.008</sup>	0.528 <sup>±0.009</sup>	0.639 <sup>±0.008</sup>	5.326 <sup>±0.2087</sup>	5.027 <sup>±0.014</sup>	6.691 <sup>±0.207</sup>	2.433 <sup>±0.059</sup>
w/o Time Modulation	0.218 <sup>±0.008</sup>	0.317 <sup>±0.009</sup>	0.420 <sup>±0.009</sup>	7.071 <sup>±0.2251</sup>	5.556 <sup>±0.017</sup>	7.263 <sup>±0.234</sup>	3.474 <sup>±0.065</sup>	w/o Time Modulation	0.372 <sup>±0.007</sup>	0.536 <sup>±0.008</sup>	0.647 <sup>±0.007</sup>	5.162 <sup>±0.2129</sup>	5.021 <sup>±0.013</sup>	7.405 <sup>±0.211</sup>	2.767 <sup>±0.062</sup>

### Observations

- **Hierarchy Matters:** Removing levels harms synergy and quality, drops R-Precision by **12 %** and worsens FID by **8 %**
- **Language Guidance is crucial** Excluding LLM cues lowers R-Precision by **9 %** and raises FID by **5 %**.
- With only 15 steps, FID is **4 %** better than InterGen while **65 %** faster; 55 steps yields peak quality.
- **Diversity & Task Conformity:** Hierarchical latent diffusion explores a broader **codebook space** for varied interactions, while LLM-guided anchors ensure **prompt consistency**.
- **SOTA Performance:** Against InterGen, R-Precision improves by **+11 %** (CORE-4D) and **+3 %** (InterHuman) and FID drops by **-5 %** and **-16 %**.

### Discussion

- We leverage hierarchical cues to accelerate alignment and generation, and our **coarse-to-fine latents reduce collisions** while preserving long-range coordination in dense scenes.
- We envision that our pipeline can support **humanoid robots** or **VR avatars**, synthesizing contact-aware trajectories from high-level goals (e.g., “hand over the wrench”).
- We plan to extend our system with **interactive motion editing** via user-adjustable handles and to integrate **scene-based cues** for more refined motion trajectories.



For more information,  
and visualizations,  
Please visit / Scan →

