

CVPR
JUNE 3-7, 2026



DENVER
COLORADO

ROMo

A Large-Scale, **R**ichly **O**rganized Dataset and
Semantic Taxonomy for Human **M**otion Generation

Jiahao Zhang^{1,2}, Joseph Liu², Young-Yoon Lee², Seonghyeon Moon²

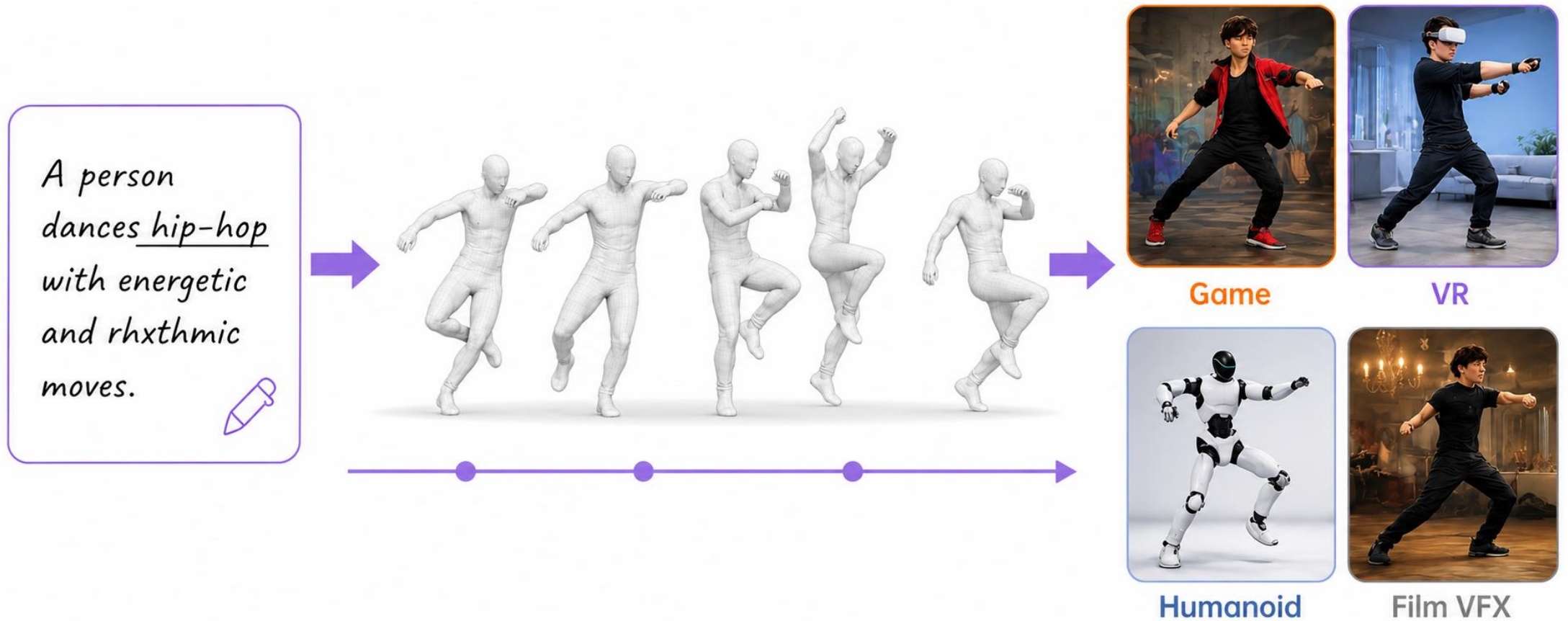
Victor Zordan², Guy Tevet³, C. Karen Liu³, Stephen Gould¹

Oren Jacob², HaomiaoJiang², Mubbasir Kapadia^{2,4}, Yizhak Ben-Shabat²

¹Australian National University ²Roblox ³Stanford University ⁴Rutgers University



Text to Motion Generation (T2M)

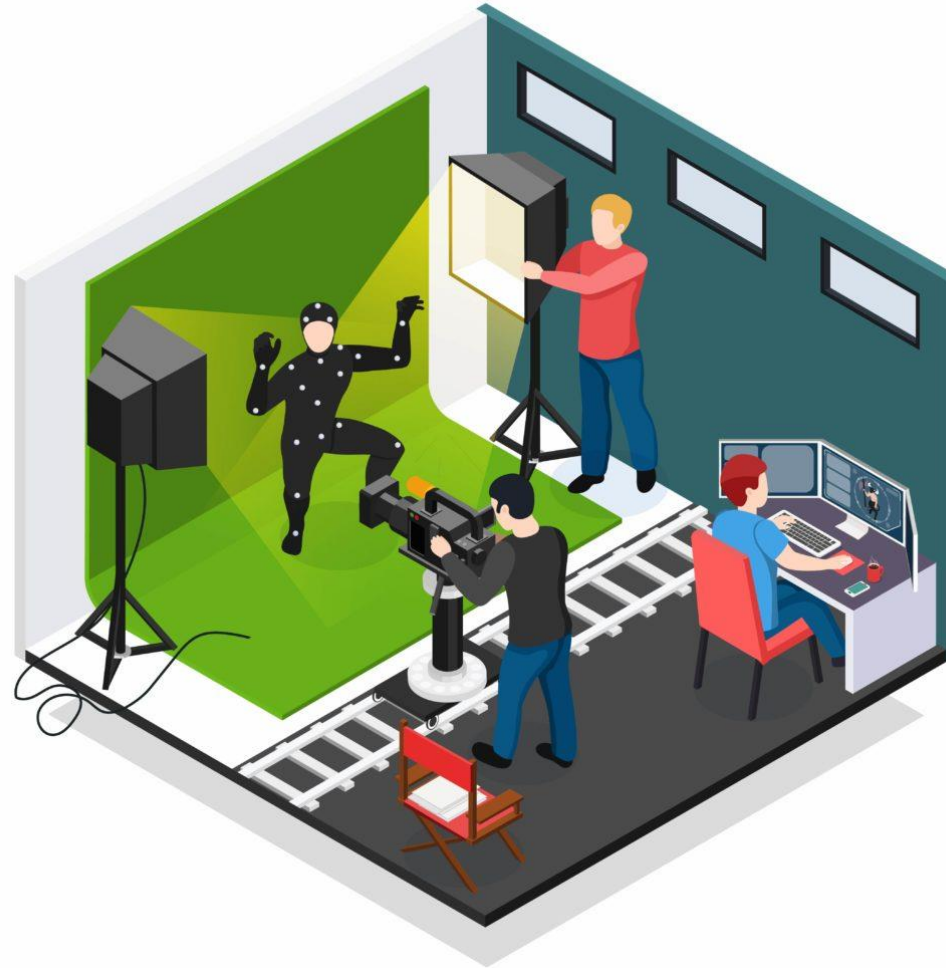


Motivation

- For images
 - ImageNet-21K (2012): **14 million** images for classification
 - LAION-5B (2022): **5 billion** text-image pairs for image generation
- What about T2M datasets?

	SEQ NUM	TEXT NUM	HOURS	MOTION	TEXT	RGB	DEPTH	BBOX	PERSON
KIT (Plappert et al., 2016)	5.7K	5.7K	11.2	B	body	✗	✗	✗	single
HumanML3D (Guo et al., 2022a)	29.2K	89K	28.6	B	body	✗	✗	✗	single
MotionX (Lin et al., 2024)	81.1K	142K	144.2	B,H,F	body	✓	✗	✗	single
MotionVerse (Zhang et al., 2024a)	320k	373k	-	B,H,F	body	✓	✗	✗	single

Motion Capture is Complicated and Expensive

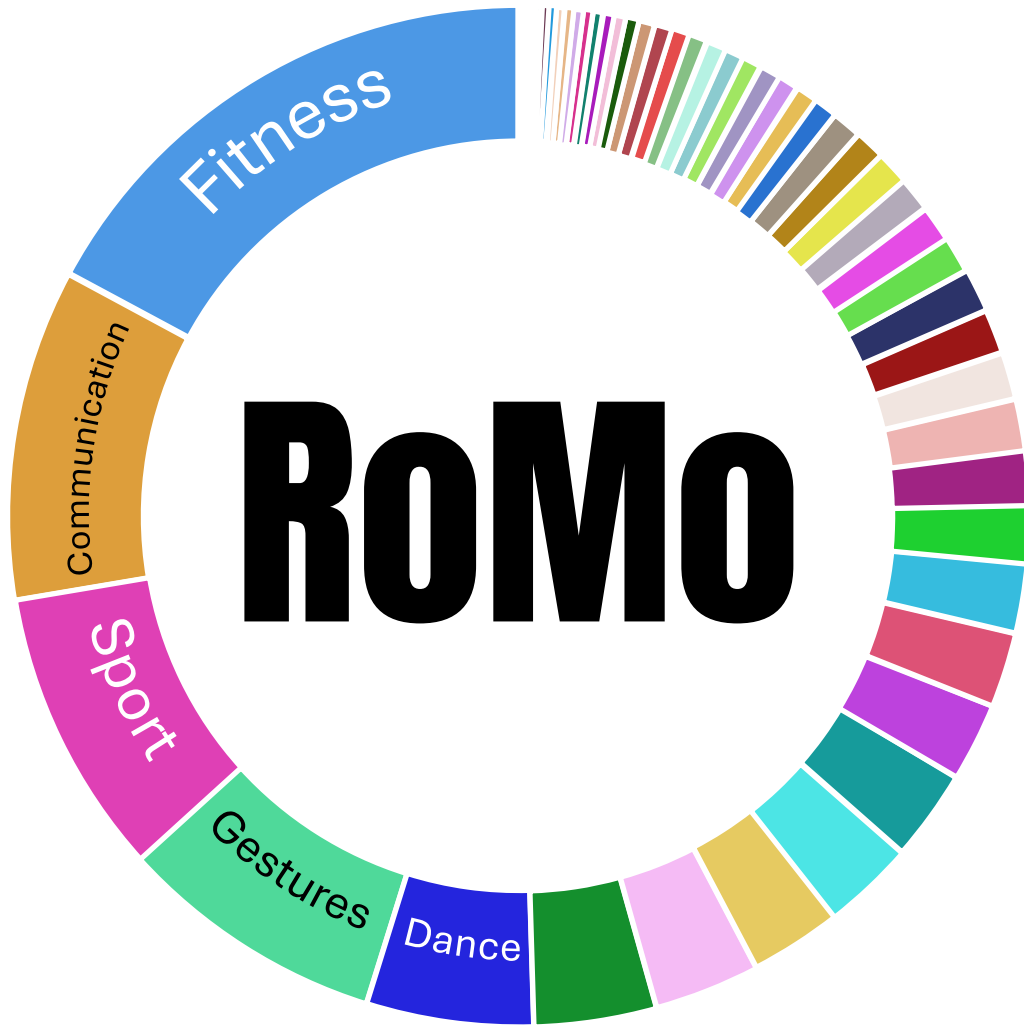


Our Solution

Same as **ImageNet** and **LAION** use the internet as the data source!

ROMO

A Large-Scale, **R**ichly **O**rganized Dataset and
Semantic Taxonomy for Human **M**otion Generation



RoMo is organized into

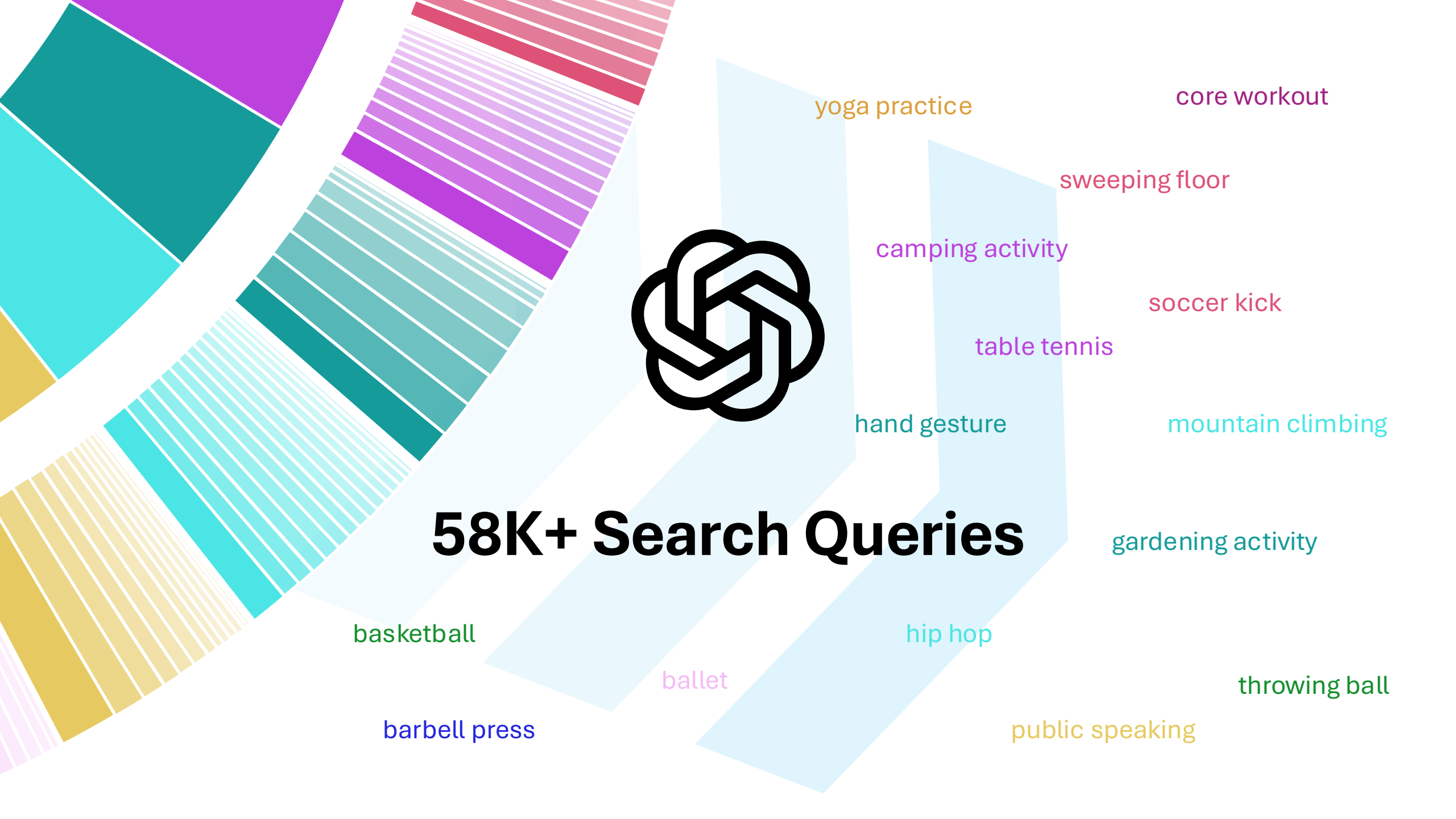
54 Categories

Further divided into
2,065 SubCategories





58K+ Search Queries



basketball

barbell press

yoga practice

core workout

hand gesture

hip hop

camping activity

sweeping floor

table tennis

soccer kick

ballet

mountain climbing

gardening activity

throwing ball

public speaking

Only **1** % remaining

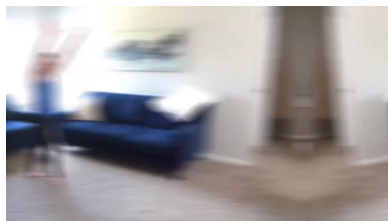
1,237+ Hours
of Human Motion Videos



Qwen3-VL

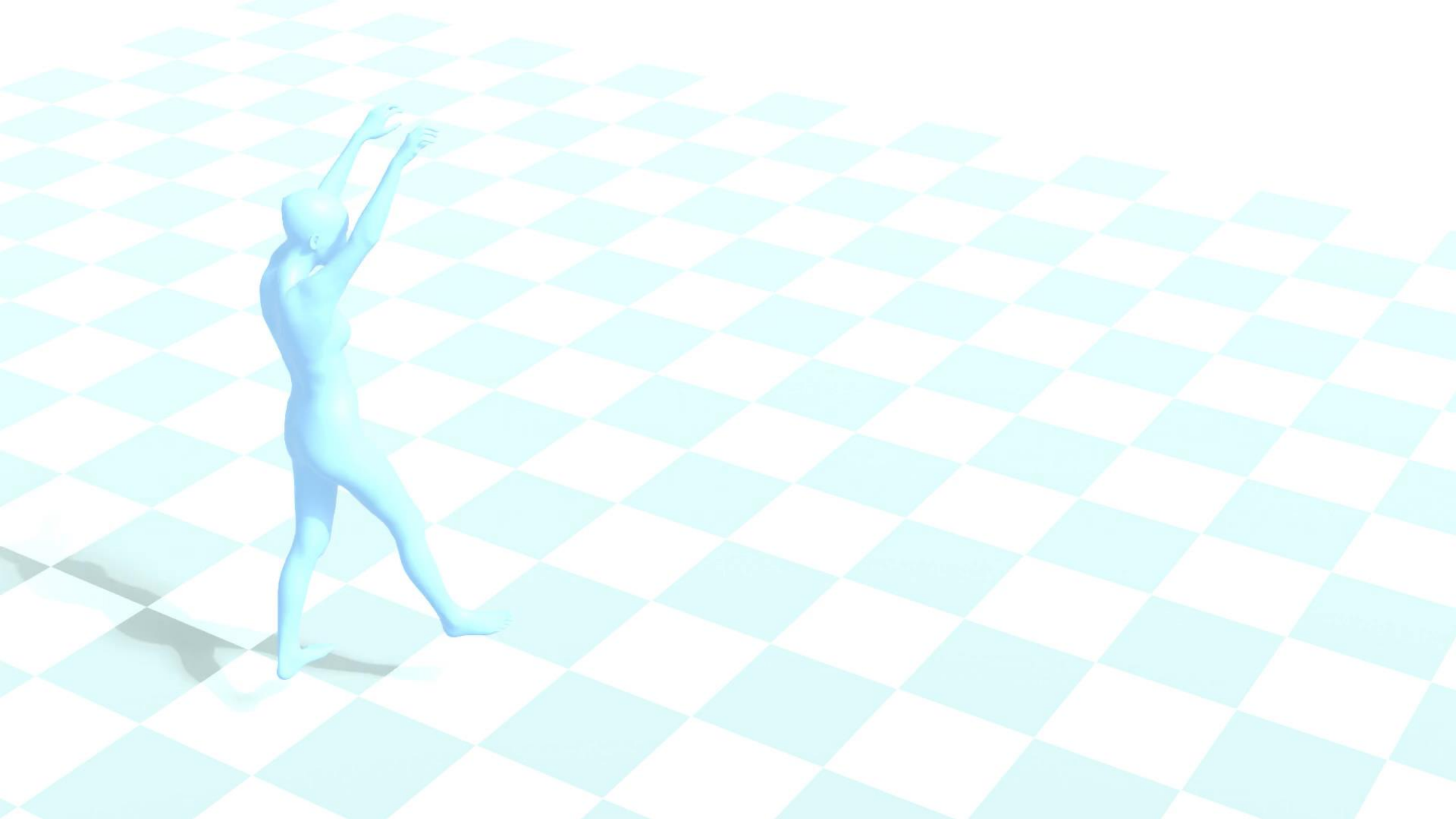


```
[{
  "start": 0,
  "end": 16,
  "action": "lift barbell"
},
{
  "start": 16,
  "end": 19.8,
  "action": "walk away"
}]
```





GVHM





Fitness → Weighing → Lift barbell overhead

“A person lifts a barbell from the ground to above his head in a smooth motion, lowering it back down repeatedly.”

Video Games → Motion Gaming → Swing

“A person in a VR headset swings a virtual sword at floating objects, hitting them and causing them to shatter.”

Sport → Bowling → Roll Bowling Ball

“A person picks up the ball, takes a few steps forward, and releases it with a smooth motion.”

Outdoors → Rafting → Push raft with paddle

“A person in a raft on the riverbank uses a paddle to push off from the shore.”

Fitness → Weighing → Lift barbell overhead

“A person lifts a barbell from the ground to above his head in a smooth motion, lowering it back down repeatedly.”



Sport → Bowling → Roll Bowling Ball

“A person picks up the ball, takes a few steps forward, and releases it with a smooth motion.”



Video Games → Motion Gaming → Swing

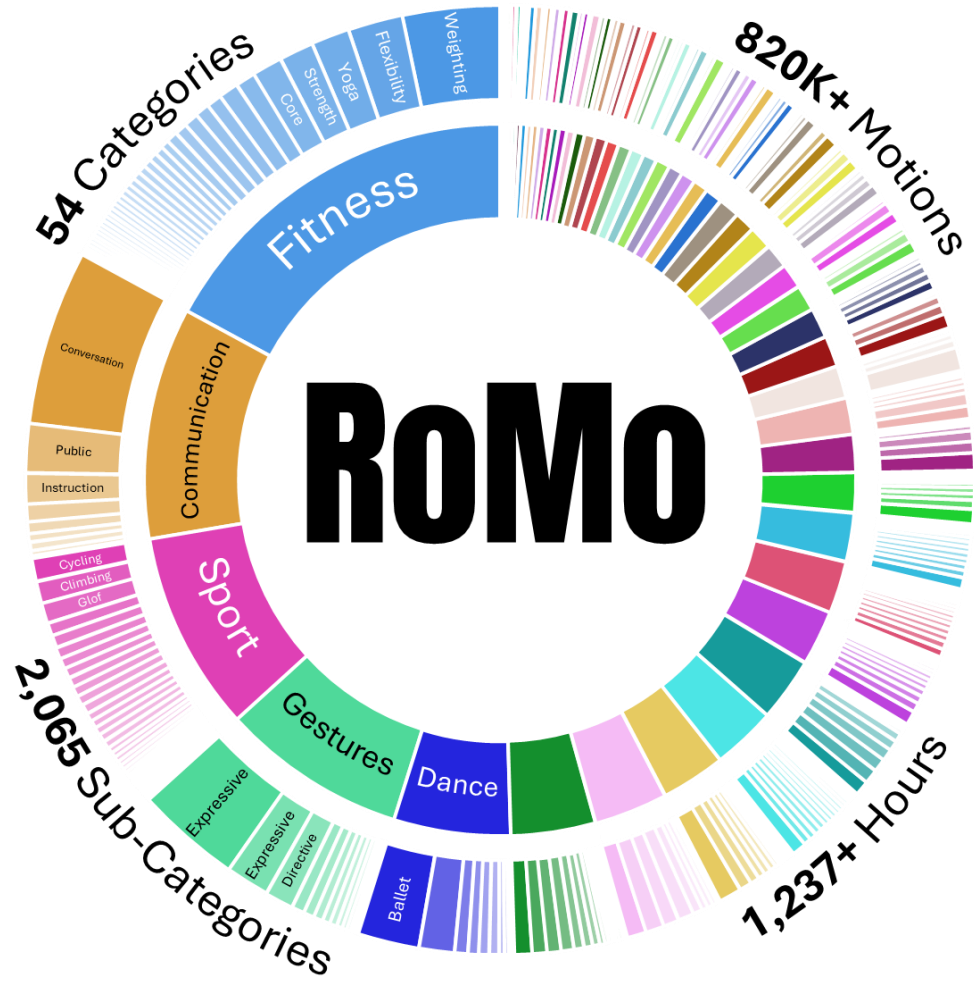
“A person in a VR headset swings a virtual sword at floating objects, hitting them and causing them to shatter.”



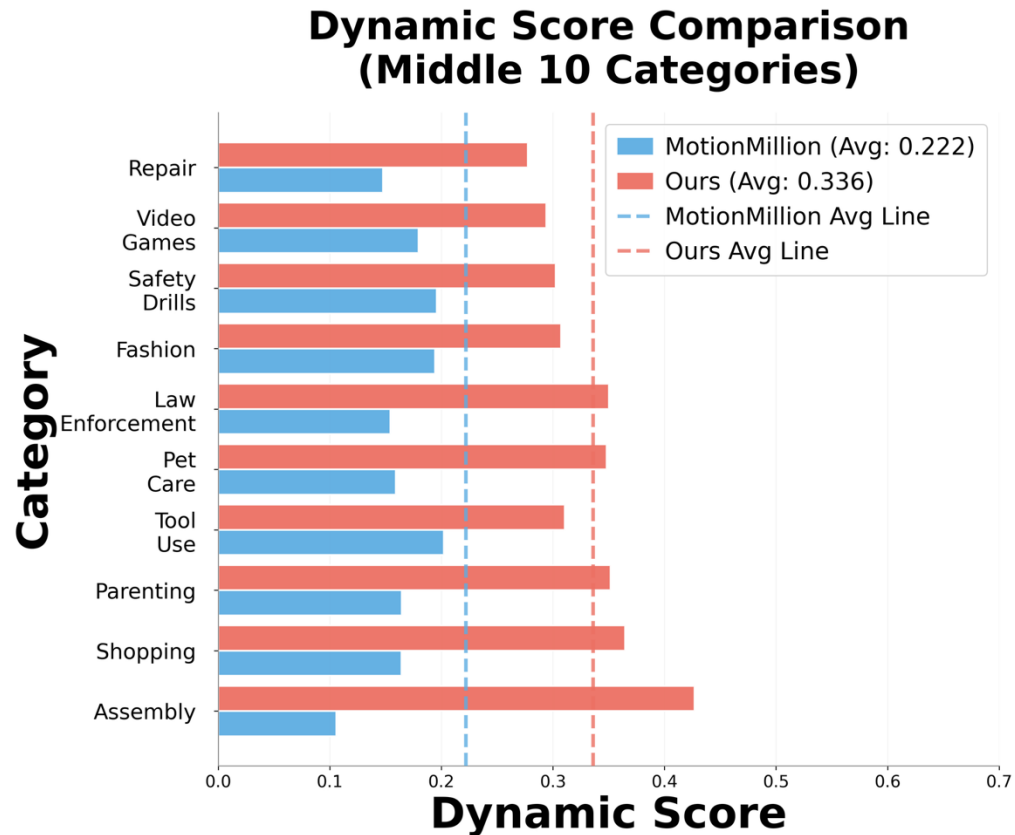
Outdoors → Rafting → Push raft with paddle

“A person in a raft on the riverbank uses a paddle to push off from the shore.”



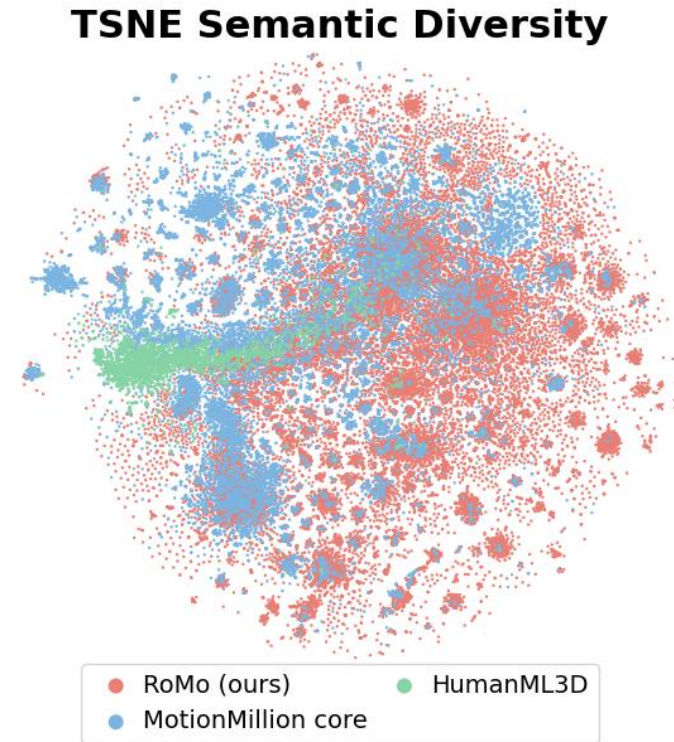
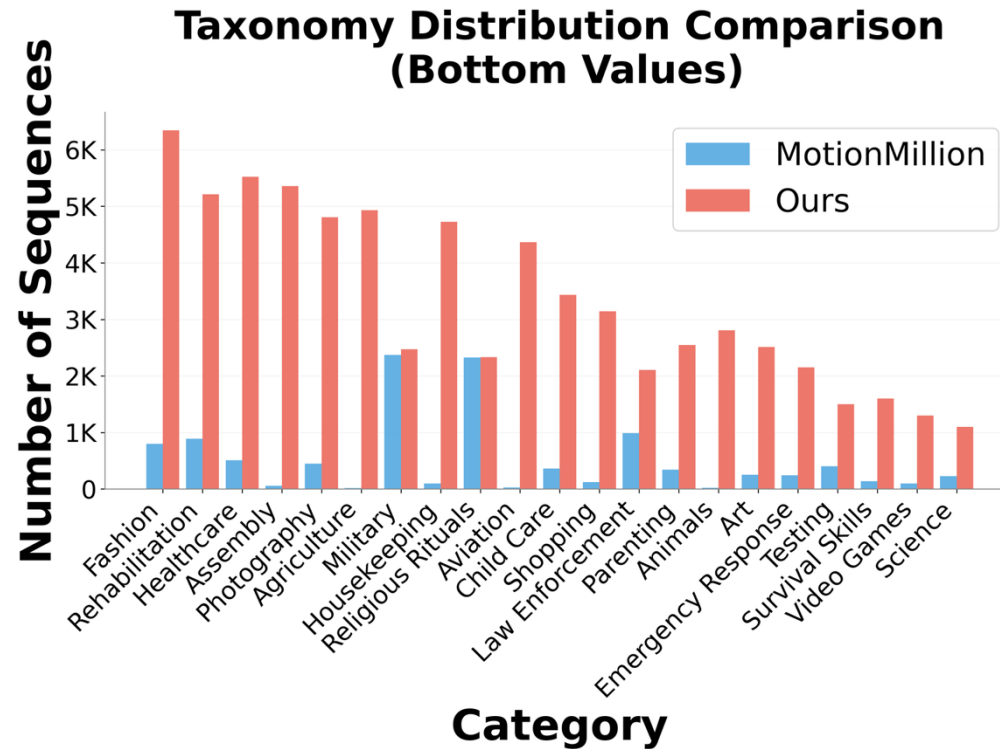


Superior Diversity, Coverage and Motion Dynamics



+41% higher dynamic score over MotionMillion.

Superior Diversity, Coverage and Motion Dynamics



+61.7% subcategory coverage over MotionMillion.

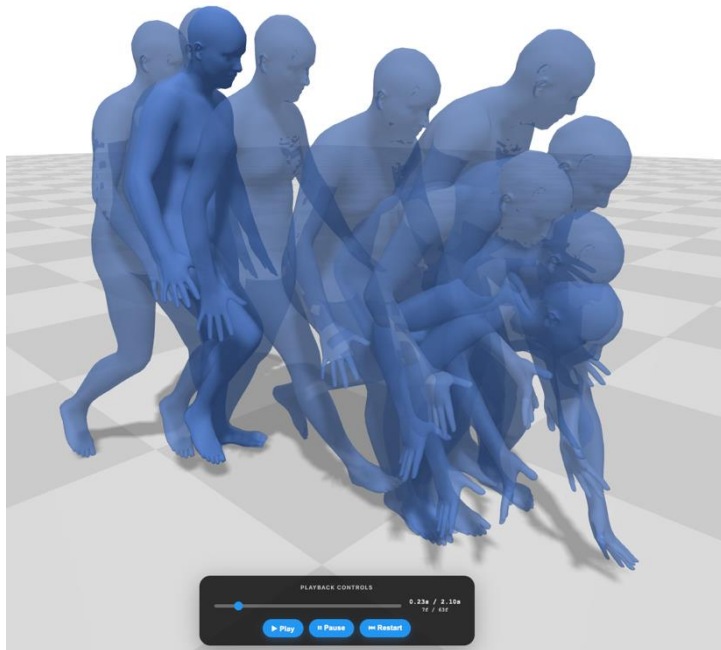
Training Modern Motion Generators on RoMo

Method	Diversity ↑	FID ↓	Matching Score ↑	Dynamic Score ↑	Ground Penetration ↓ ($\times 10^{-5}$)	Foot Skating ↓ ($\times 10^{-3}$)	Floating ↓ ($\times 10^{-2}$)
MDM	27.67	20.63	12.06	0.2138	0.0	1.70	1.67
MMGPT	16.68	12.80	22.08	0.3268	3.55	92.0	0.0311

Open-Source Motion Toolbox

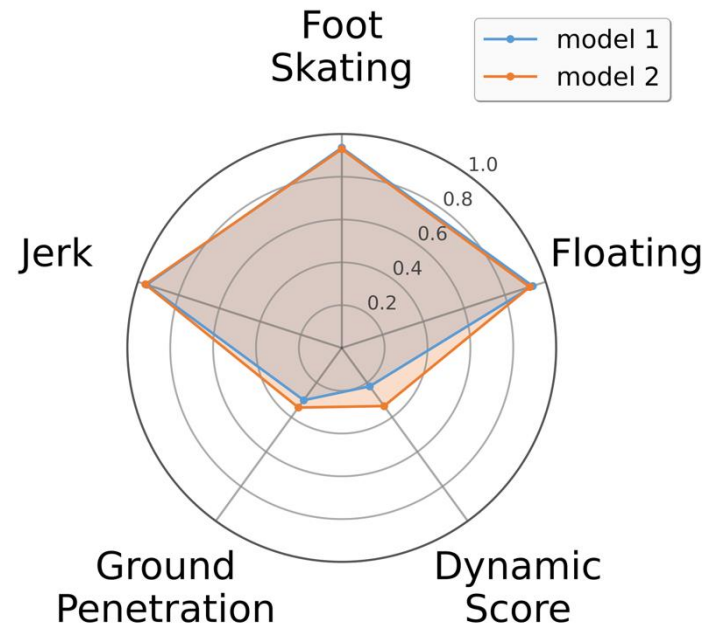
Interactive Visualization

Inspect and compare motion sequences directly in the browser.



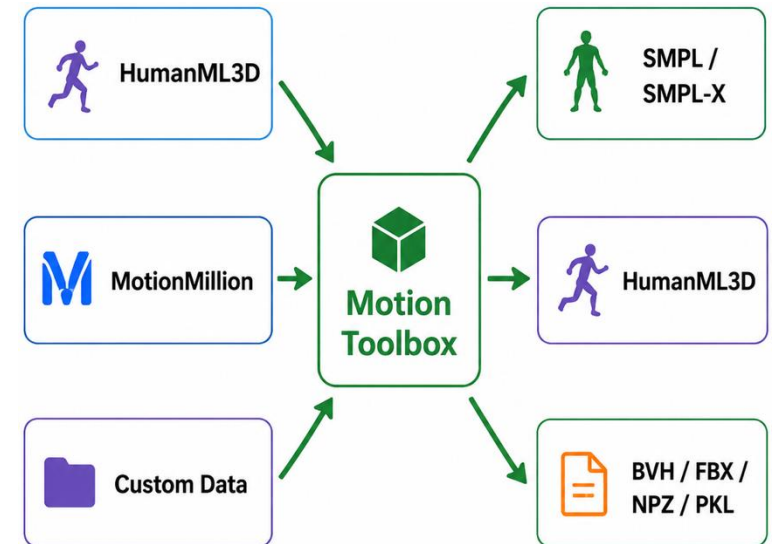
Unified Benchmarking

Reproducible evaluation tools for large-scale human motion research.



Unified Conversion

Convert between popular motion formats and research datasets.





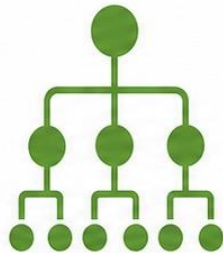
Thanks!



820K | **1,237**
Clips | Hours

Large-Scale Motion Dataset

Rich captions and high-quality filtered motion.



54 | **2065**
Category | Subcategory

Taxonomy-Aware Curation

Adaptive filtering for quality and diversity.



**Motion
Toolbox**

Open Source Motion Toolbox

Standardized evaluation and motion analysis.