

Motion Seed Experiment v2 Report

Date : 2026-02-14 Total trials : 700 (140 seeds x 5 LLMs) Runtime : ~60 min
~60 min Prompt version : v2 (enriched: control law + weight semantics + few-shot + criteria hints + discretization)

Experiment Design

28 motion concepts x 5 languages (en/de/zh/fr/fi) = 140 seeds, prompted to 5 LLMs. Each LLM generates 6 neural network weights snapped to a 9-point grid {-1.0, -0.7, -0.4, -0.1, 0, 0.1, 0.4, 0.7, 1.0}. Weights drive a headless PyBullet simulation (4000 steps @ 240 Hz) of a 3-link robot, and the resulting gait is evaluated against behavioral criteria.

Prompt Improvements over v1

The v2 prompt was designed by consulting three frontier LLMs (GPT-4.1-mini, DeepSeek, GPT-5.2) for improvement strategies:

Feature	v1	v2	Source
Weight semantics	None	Full sensor→motor mapping with functional labels	DeepSeek + GPT-4.1-mini
Control law	None	`m3 = tanh(w03*s0 + w13*s1 + w23*s2)`	GPT-5.2
Few-shot examples	None	3 contrasting verified examples with measured outcomes	All three
Design priors	None	Structural templates (symmetry, cross-coupling, directional bias)	GPT-5.2
Behavioral hints	None	Qualitative per-concept (not numeric thresholds)	DeepSeek
Chain-of-thought	Forbidden ("No explanation")	In 1-2 sentences, note the motor pattern"	DeepSeek
Weight discretization	Continuous [-1, 1]	Snapped to 9-point grid	GPT-5.2
Token budget	200 / 1000	500 / 2000	Empirical

Models

Model	Type	Parse Rate	Concept Match Rate
gpt-4.1-mini	OpenAI API	140/140 (100%)	33/140 (24%)
qwen3-coder:30b	Ollama (local)	140/140 (100%)	22/140 (16%)
gpt-oss:20b	Ollama (local)	133/140 (95%)	22/131 (17%)
llama3.1:latest	Ollama (local)	136/140 (97%)	21/136 (15%)
deepseek-r1:8b	Ollama (local)	118/140 (84%)	22/114 (19%)

28 Motion Concepts — Match Grid

Concept	qwen3-coder	deepseek-r1	llama3.1	gpt-oss	gpt-4.1-mini
freeze	5/5 (100%)	5/5 (100%)	2/5 (40%)	5/5 (100%)	5/5 (100%)
drag	3/5 (60%)	1/4 (25%)	2/5 (40%)	3/5 (60%)	1/5 (20%)
crawl	0/5	1/5 (20%)	1/5 (20%)	2/5 (40%)	4/5 (80%)
retreat	2/5 (40%)	2/4 (50%)	0/5	1/4 (25%)	4/5 (80%)
hop	1/5 (20%)	2/4 (50%)	1/5 (20%)	3/5 (60%)	3/5 (60%)
patrol	1/5 (20%)	1/5 (20%)	3/5 (60%)	1/5 (20%)	1/5 (20%)
sway	0/5	2/5 (40%)	3/4 (75%)	1/5 (20%)	0/5
wobble	1/5 (20%)	0/4	2/5 (40%)	2/5 (40%)	0/5

zigzag	0/5	2/5 (40%)	2/4 (50%)	0/5	1/5 (20%)
gallop	1/5 (20%)	1/5 (20%)	0/5	1/5 (20%)	1/5 (20%)
stagger	1/5 (20%)	2/3 (67%)	1/5 (20%)	0/5	0/5
dash	1/5 (20%)	1/5 (20%)	1/4 (25%)	0/5	0/5
sprint	2/5 (40%)	0/5	0/5	1/5 (20%)	0/5
stomp	2/5 (40%)	0/5	0/5	1/4 (25%)	0/5
march	1/5 (20%)	0/4	0/5	1/4 (25%)	0/5
drift	0/5	0/3	0/5	0/5	2/5 (40%)
tiptoe	0/5	0/5	0/5	0/5	2/5 (40%)
charge	1/5 (20%)	1/4 (25%)	0/5	0/5	0/5
turn_left	2/5 (40%)	0/4	1/5 (20%)	1/5 (20%)	0/5
scurry	0/5	1/4 (25%)	1/5 (20%)	0/5	0/5
plod	0/5	0/4	1/5 (20%)	0/5	0/5
rock	0/5	0/2	0/5	1/5 (20%)	0/5
slide	0/5	0/4	1/5 (20%)	0/5	0/5
headstand	0/5	1/2 (50%)	0/5	0/5	0/5
turn_right	0/5	0/4	0/5	0/5	1/5 (20%)
pivot	0/5	0/4	1/5 (20%)	0/4	0/5
twirl	0/5	0/5	1/5 (20%)	0/5	0/5
circle	0/5	0/4	0/5	0/5	0/5

Key Findings

1. Parse rate dramatically improved

v2 achieved **97% parse success** across all models (680/700), vs v1's initial 56% (before fixes). The key factors:

- Weight discretization instruction ("Choose each weight from [grid]") anchored model output
- Brief CoT instruction ("In 1-2 sentences") prevented verbose reasoning from consuming the token budget
- Increased token budget (500/2000 vs 200/1000) gave reasoning models room

2. "Freeze" is the new "stand still" — 100% on 4/5 models

All 5 LLMs across all 5 languages independently discovered that all-zero weights = no motion for "freeze." This replicates v1's "stand still" finding and extends it: the enriched prompt with the control law equation ($m3 = \tanh(w03*s0 + \dots)$) makes it trivially obvious that zero weights → zero motor commands → no motion. Even llama3.1 matched 2/5 (vs 0/5 for stand_still in v1).

3. Retreat >> backward_walk — reframing helps directional control

"backward_walk" in v1 scored 15% (3/20). "retreat" in v2 scored **9/23 (39%)** with matches across 4 languages and all 5 models. The qualitative hint "backward movement, moving away from the forward direction" was more effective than the implied directional instruction. gpt-4.1-mini matched 4/5 languages for

retreat, suggesting it genuinely associates "retreat" with negative-DX weight patterns.

4. Drag is universally understood

"drag" achieved 10/24 (42%) – the highest match rate for a non-trivial motion concept. Chinese "拖行" was particularly effective: all 5 models matched. The hint "slow effortful movement, lots of energy for little progress" gave models a clear behavioral target.

5. GPT-5.2's "doorbells for ballet" prediction confirmed

Circle (0/24), twirl (1/25), and slide (1/24) scored near zero. These require *timing and continuous control*, not touch-triggered reflexes. The robot's reflex-only controller physically cannot produce smooth circular paths or gliding. This validates GPT-5.2's insight that the controller's expressivity – not prompting – is the bottleneck for these concepts.

6. Weight discretization reduced but didn't eliminate collapse

qwen3-coder still shows some weight-vector collapse (e.g., the vector [-0.22, +0.47, +3.2...] appears repeatedly), but less severely than v1. The 9-point grid constrains the space from infinite floats to $9^6 = 531,441$ possible vectors, making collapse more detectable and diversity more achievable.

7. Headstand: exactly 1 match

Only deepseek-r1:8b matched headstand, via Finnish "seisoa päällään" (literally "stand on one's head"). The resulting weights [+0.7, -0.7, +0.1, +1.0, -0.7, +0.4] produced torso_duty > 0.6, meaning the robot's body was on the ground (inverted) most of the simulation. This is the Ubik gait pattern – strong opposing torso weights flip the robot.

8. Turn left/right: asymmetric success

Turn_left scored 4/24 (17%) with matches in German, Finnish, and English. Turn_right scored only 1/24 (4%) – a single English match from gpt-4.1-mini. This asymmetry is surprising and may reflect training data biases or accidental correlation between "left" and negative-yaw weight patterns.

9. Model specializations emerged

Model	Strength	Weakness
gpt-4.1-mini	crawl (80%), retreat (80%), freeze (100%), drift, tiptoe	sway, wobble, stagger
deepseek-r1:8b	stagger (67%), headstand (50%), hop (50%)	Low parse rate (84%), no matches for many concepts
llama3.1:latest	sway (75%), patrol (60%), zigzag (50%)	freeze (only 40%), charge, retreat
gpt-oss:20b	freeze (100%), drag (60%), hop (60%)	sway, stagger, dash
qwen3-coder:30b	freeze (100%), drag (60%), sprint (40%)	sway, crawl, plod

10. The combined dictionary has 58 concepts and 365 entries

Merging v1 (12 original concepts + 40 English extras) with v2 (28 new concepts), the motion gait dictionary now covers 58 distinct motion concepts with 365 viable weight-vector entries, each citing the descriptive word, model, and language.

Comparison: v1 vs v2

Metric	v1	v2
Prompt length	350 chars	1,646 chars
Parse rate (all models)	83% (after fix)	97%
Parse rate (deepseek-r1)	17%	84%
Concepts tested	12 core + 40 extras	28 core
Unique concepts matched	12/12 (100%)	24/28 (86%)
Unmatched concepts	(none)	circle, slide (partial), rock (partial), plod (partial)
gpt-4.1-mini match rate	43% (core)	24% (but harder concepts)
Best single-concept rate	stumble 74%	freeze 86%
Worst single-concept rate	forward_walk 9%	circle 0%

The lower overall match rate in v2 reflects harder concepts (circle, twirl, headstand), not worse prompting. For comparable concepts, v2 generally improved: retreat (39% vs backward_walk 15%), crawl (32% vs not tested), hop (40% vs bounce 35%).

Technical Notes

- **Weight discretization** : All weights snapped to nearest value in {-1.0, -0.7, -0.4, -0.1, 0, 0.1, 0.4, 0.7, 1.0} after parsing. Models frequently generated exact grid values voluntarily.
- **CoT token management** : The instruction "In 1-2 sentences, note the motor pattern" + "Keep reasoning SHORT" successfully constrained output to ~100-200 chars of reasoning before JSON, vs v1's initial failure where models consumed all tokens on reasoning.
- **Prompt consultations** : Three frontier LLMs (GPT-4.1-mini via API, DeepSeek via web chat, GPT-5.2 via web chat) were consulted for prompt design. Key insight from GPT-5.2: "Some words want timing, phase, and smooth continuous control. If your controller can't make a rhythm without contact events, you're asking it to compose ballet using only doorbells."

Files

- `motion_seed_experiment_v2.py` – Experiment script (28 concepts, enriched prompt)
- `build_motion_dictionary.py` – Dictionary builder (merges v1 + v2)
- `artifacts/motion_seed_experiment_v2.json` – Raw results (700 trials)

- `artifacts/motion_gait_dictionary_v2.json`
 - Combined dictionary (58 concepts, 365 entries)