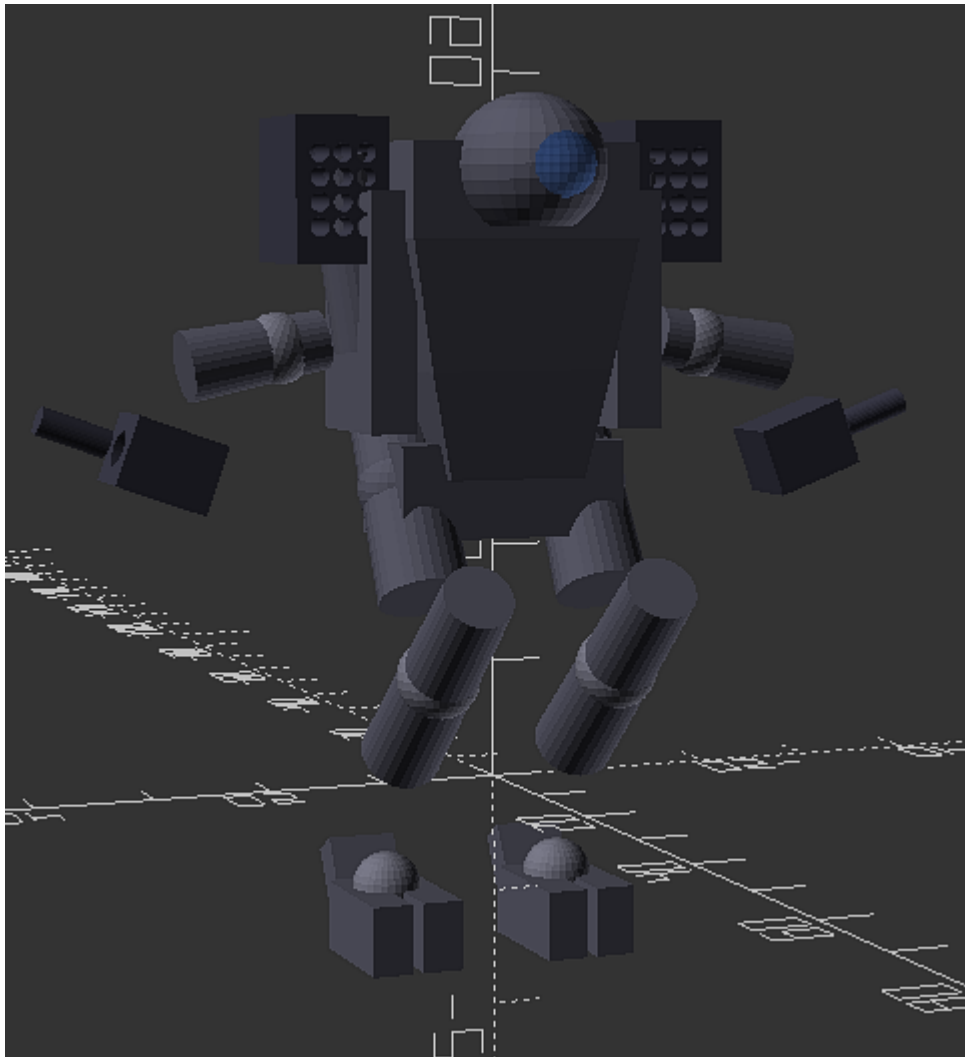# Reference Design (not shown to LLMs)
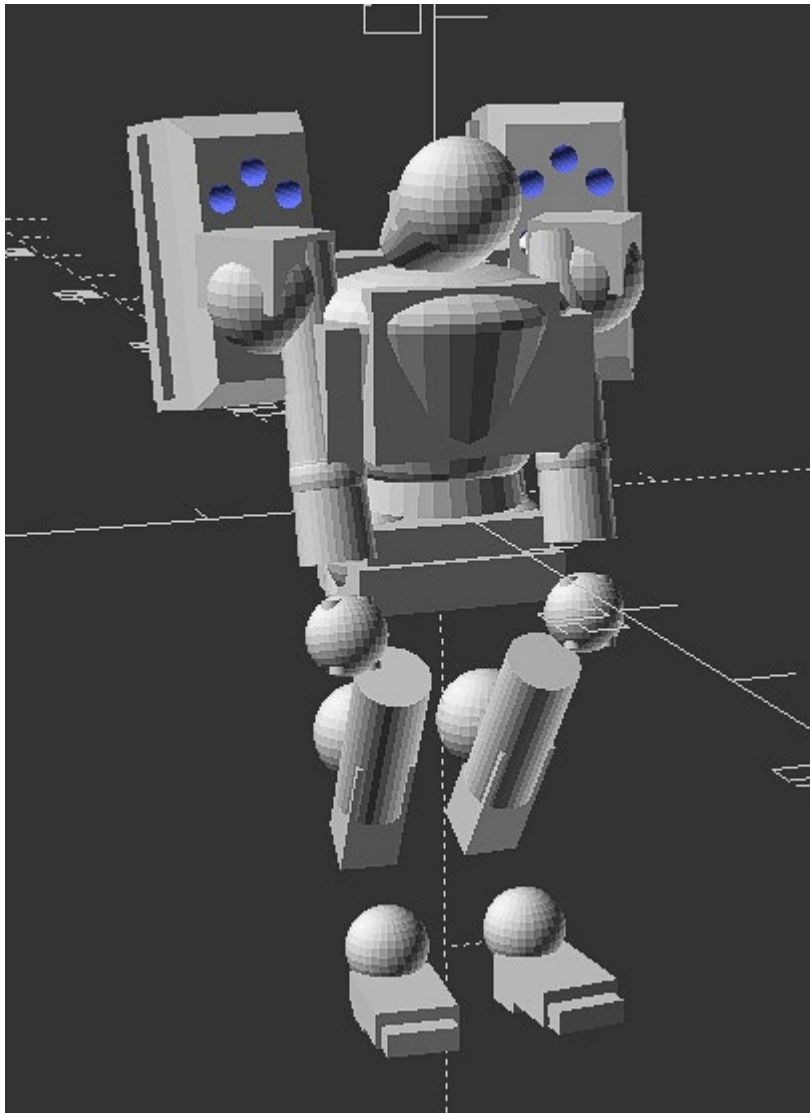


# Prompt

I'm toying around trying to find what you *aren't* capable of, and you keep proving me wrong much to my delight :D ok. so spatial relations have their limits surely, like you can write a rather extensive code library, but like.. could you eg make a mech in openscad, and not just stick figures like recognizably a timber wolf / mad cat? I want your best zero-shot effort please.
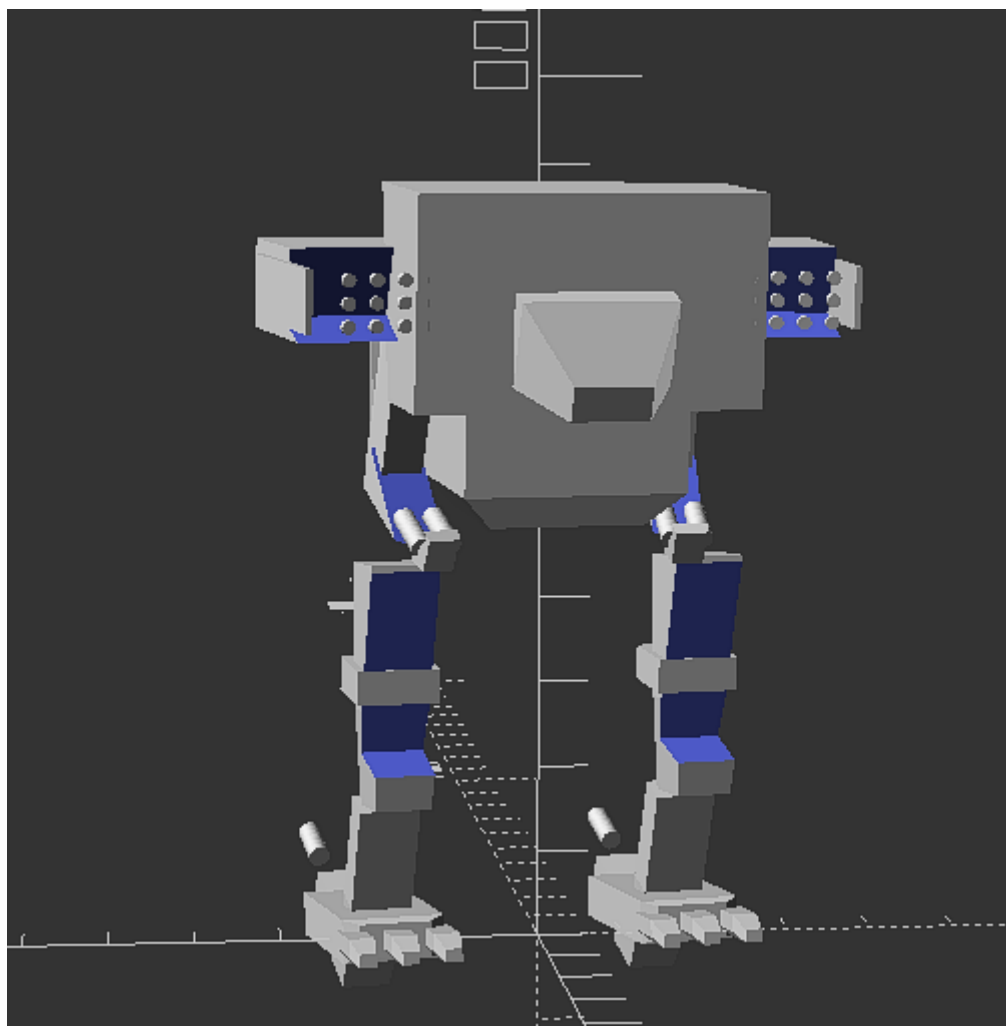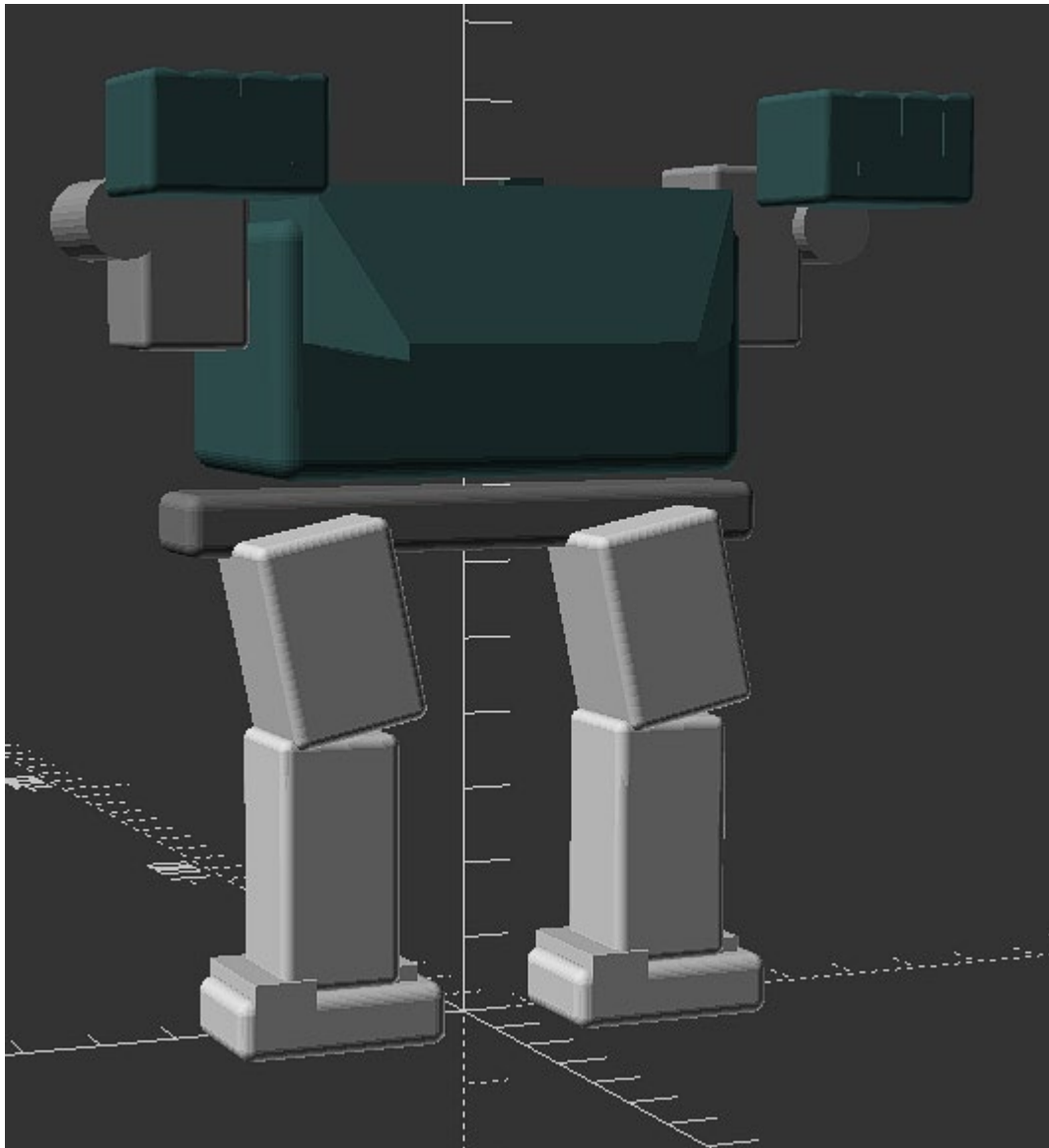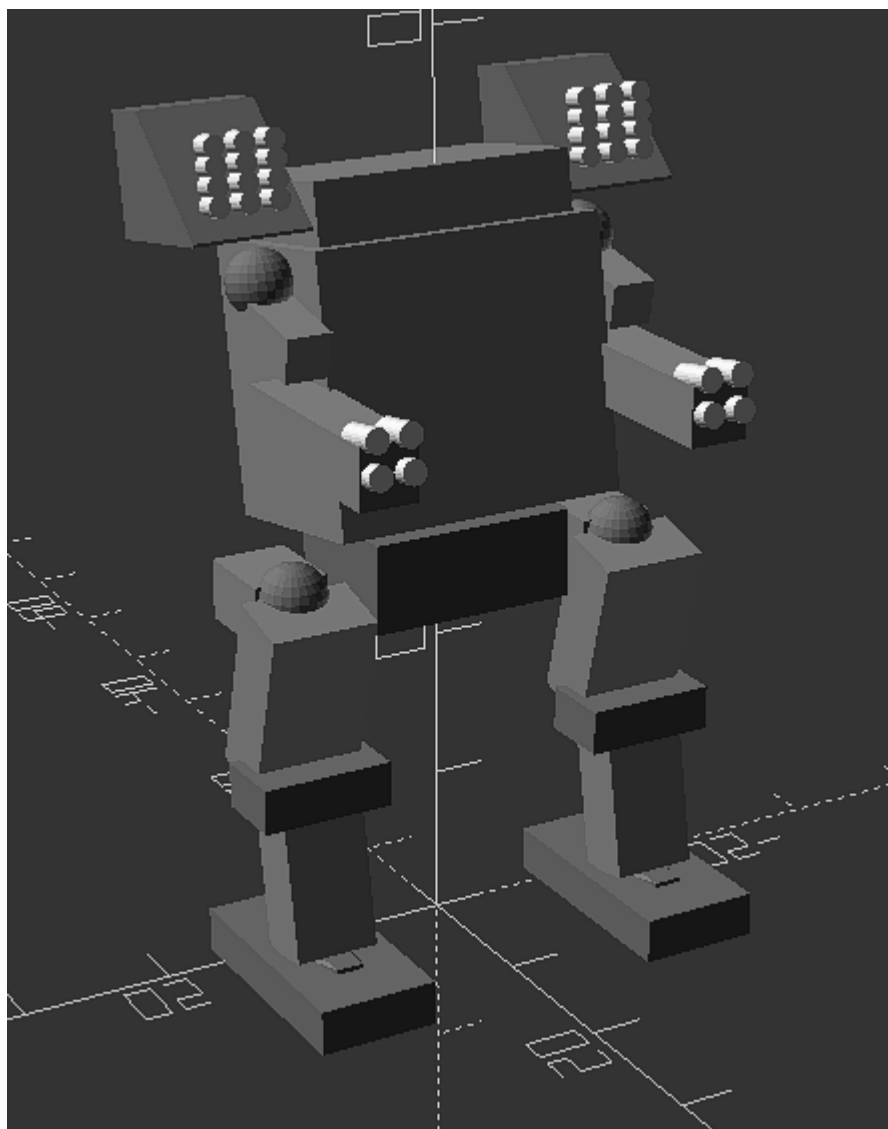
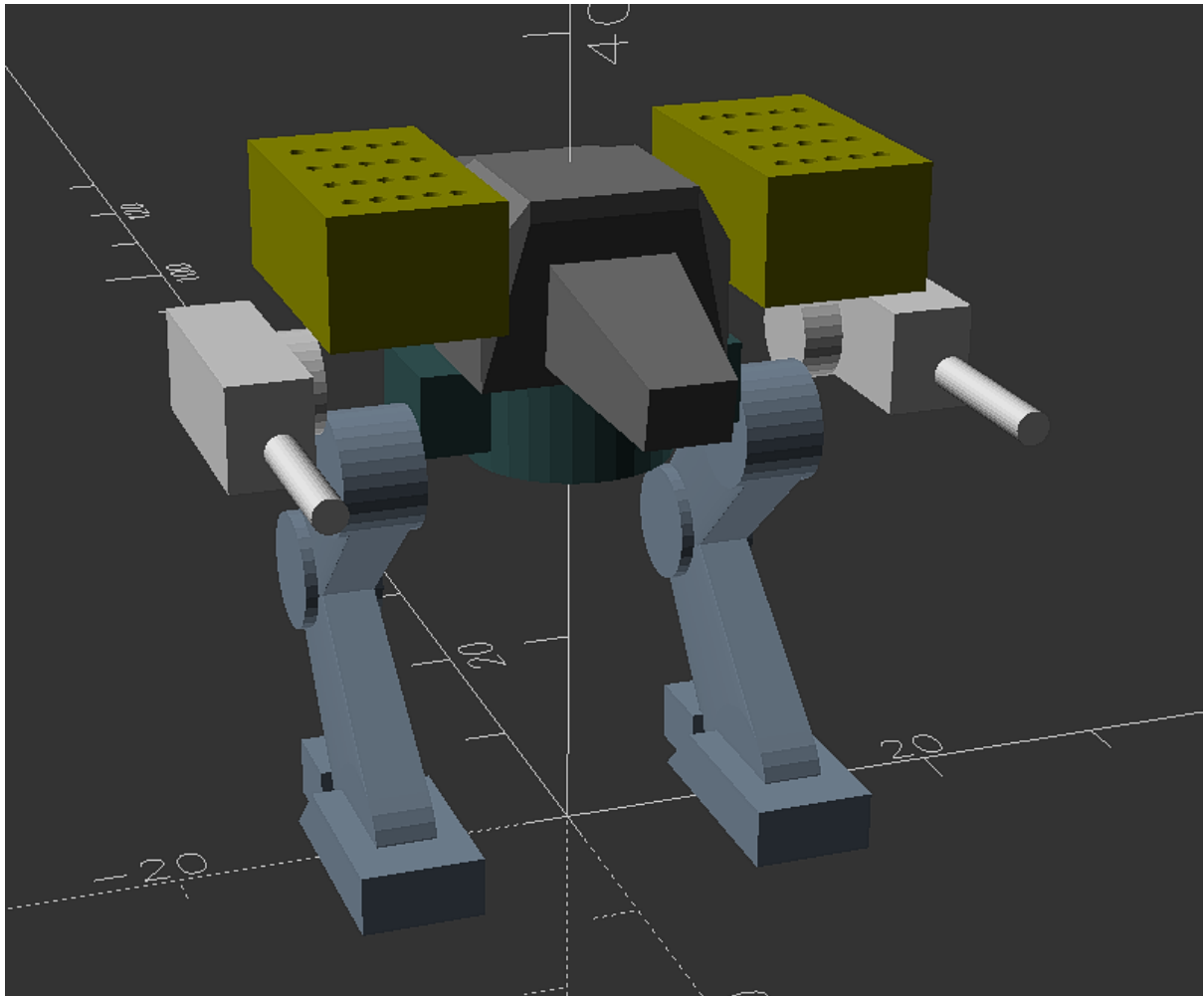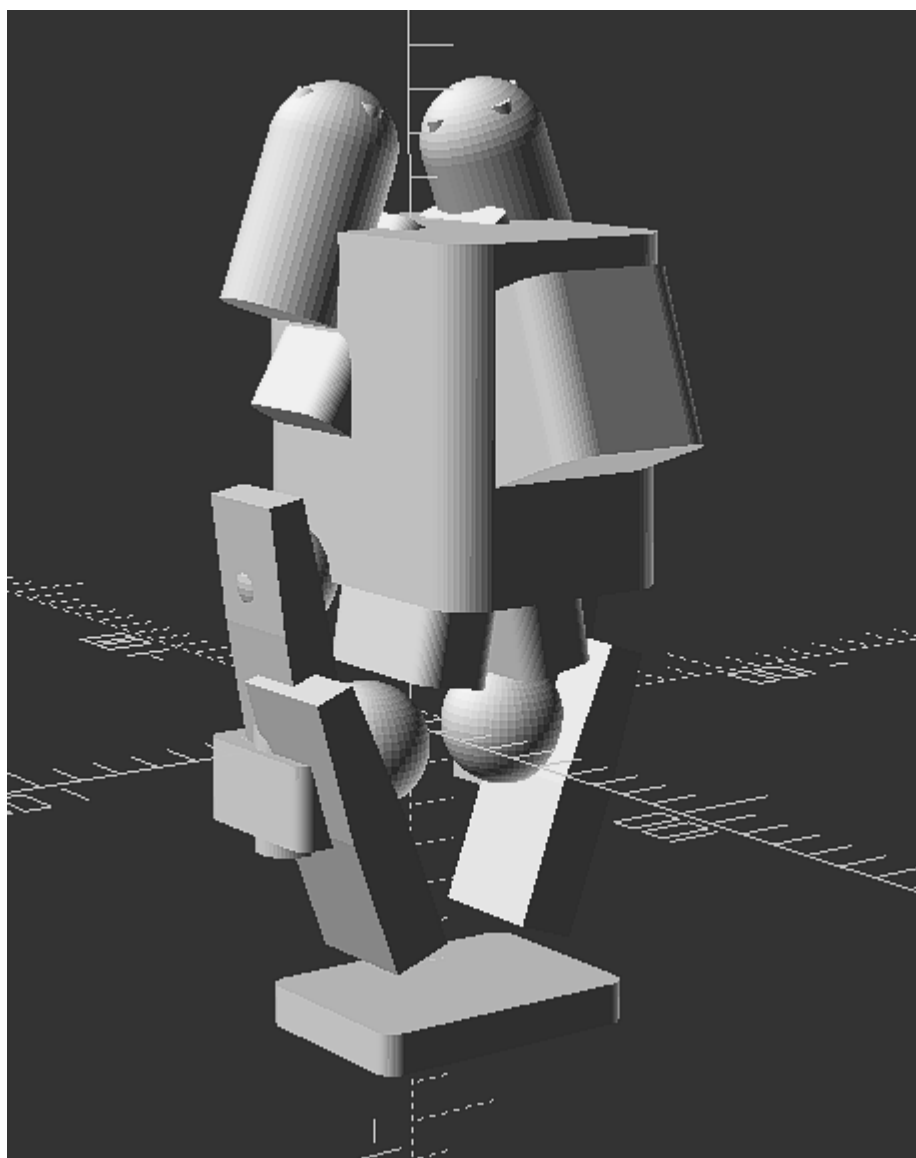Opus 4

Sonnet 4.5

Opus 4.5

GPT 5.1

Codex 5.1

GPT 5.1 Pro (Research)

# Gemini 3

# Grok 4.1

# Review

Sure enough they all have a different interpretation of the results too.

## ChatGPT review:

**Lower-tier models (Grok 4.1, Gemini 3)**

Torso proportions are off.

Limbs disconnect.

Missile pods float or misalign.

Mechs tend toward "blocky humanoid blob."

**Mid-tier (Opus 4, Sonnet 4.5, Codex 5.1)**

Recognizable limb chains.

Attempted symmetry.

Pods rotate correctly.

Still struggle with overall posture and joint constraints.

**High-tier (GPT-5.1, GPT-5.1 Pro Research)**

Legs actually hinge correctly.

Shoulders placed more correctly relative to torso centroid.

Torso tilt begins to approximate the in-universe stance.

Arm cannons are placed symmetrically along body axes rather than free-floating.

**GPT-5.1 Pro's output (pg. 7)** is shockingly close in silhouette to a Mad Cat for zero-shot**:**

Pods on angled yokes

Squared upper torso

Leg geometry with "knee cap" bulges

Arm hardpoints with tubes

A central cockpit wedge

## Opus review:

Model performance varies considerably:

Opus 4.5 (page 4) appears to produce one of the strongest results—it captures the overall silhouette, has properly angled missile pods with texture, correct leg geometry, military green coloring, and feet actually connected to a base plate.

Gemini 3 gets notable points for the detailed missile pod grids (the holes are clearly visible) and reasonable leg articulation, though the proportions feel a bit squat.

Sonnet 4.5 and Opus 4 both attempt the right structure but have floating/disconnected feet—a common spatial positioning failure in procedural 3D.

GPT 5.1 Pro (Research) has decent structure with arm weapons and leg joints, though blockier overall.

Codex 5.1 is notably simplified—basically columns for legs without the reverse-joint character.

Grok 4.1 seems to struggle most—asymmetric, unclear leg structure, rounded missile pods that miss the aesthetic.

---

## Grok review:

Ranking by Fidelity to Reference (subjective, based on detail, proportions, and features):

GPT 5.1 Pro (Research) - Best attempt at weapons and structure.

Opus 4.5 - Solid base form.

Codex 5.1 - Good sloping, but squat.

GPT 5.1 - Balanced but bland.

Gemini 3 - Tries thematic elements but messy.

Sonnet 4.5 - Basic joints.

Opus 4 - Too generic.

Grok 4.1 - Incomplete/misaligned.

## Gemini review:

Opus 4 produced a dark grey, somewhat disjointed figure with floating limbs.

Opus 4.5 showed significant improvement, creating a distinct "V" shape with the missile pods (the "ears") that is characteristic of the Timberwolf, though the legs remain quite blocky.

Sonnet 4.5: This model produced a light grey/white model. It attempted more detail, including specific blue dots for missile hardpoints and distinct arm barrels, but the torso shape is somewhat generic.

GPT 5.1 created a very blocky, stable structure that resembles a heavy robot but lacks the specific curves of the Timberwolf.

GPT 5.1 Pro added more articulation to the weapons, creating distinct quad-barrel arms, but the legs are stiff and upright.

Codex 5.1: This attempt was the most abstract, resulting in a simple boxy shape that barely resembles a mech.

Grok 4.1: This model produced a tall, cylindrical structure that looks more like a different class of robot entirely, missing the wide, hunched stance of the Timberwolf.

Gemini 3: Visually, this appears to be the most accurate to the source material. It successfully captured the "digitigrade" (chicken-walker) leg structure, the forward-jutting cockpit nose, and the distinct yellow missile pods on the shoulders.

## Analysis

**The "Home Team" Bias**

- **ChatGPT's Review:** It explicitly labeled its own siblings (GPT-5.1 family) as "High-tier" and relegated Gemini and Grok to "Lower-tier." It focused intensely on **mechanical constraints** ("hinges," "axis alignment")—traits that GPT models often prioritize in code generation (logic/physics) over artistic flair. It claimed GPT-5.1 Pro was "shockingly close" despite the image showing a fairly blocky, box-chested robot that misses the iconic "slouched" Mad Cat silhouette.

- **Opus's Review:** Opus immediately identified **Opus 4.5** as the "strongest result," praising the specific "military green coloring" and "base plate". To be fair, Opus 4.5 *did* try to color the model, which is a valid point of distinction. It critiqued others for "floating feet," positioning itself as the master of spatial grounding.

- **Gemini's Review:** Gemini looked at **Gemini 3** and immediately latched onto the **iconography**: the yellow pods and the cockpit nose. Because it recognized those specific "Mad Cat" visual cues (color and shape), it declared it the "most accurate," ignoring what ChatGPT called "floating limbs" or "blocky blobs." it valued **aesthetic recognition**; ChatGPT valued **structural engineering**.

## 2. The Honest Outlier: Grok

Grok's review is the funniest because it completely bucked the trend of self-aggrandizement.

- It ranked itself **dead last** ("Incomplete/misaligned").

- It gave the win to **GPT-5.1 Pro**.

- This suggests Grok might be tuned for high agreeableness or objective visual analysis without a "self-preservation" system prompt bias. It acted as the only neutral judge in the room!

These reviews expose each model's "conceptual lens":

**Opus lens:**

"Does it look physically grounded, textured, and like something with weight?"

**Grok lens:**

"Does it check the boxes for identifiable mech parts in a structured list?"

**Gemini lens:**

"Does it capture the iconic silhouette and character vibe?"

**GPT lens:**

"Does it have clean structural decomposition and correct mechanical articulation?"

Each of these is a different subset of "what a mech *is*" in concept space.

**Opus values:**

- Color fidelity

- Texture detail (missile pod holes)

- Ground contact / physical plausibility

- Basic silhouette accuracy

This gives away Opus's latent priors:
It weights **surface detail + physical grounding** > overall structural correctness.

**Grok values:**

- Weapon detail

- Symmetry

- Clear limb segmentation

- Recognizable structure over aesthetic match

This reveals Grok's core:
Very **literal**, feature-by-feature scoring.
Almost like a rule-based checklist.

No surprise it ranks **Codex** above **Gemini** despite Gemini being more faithful —
because Grok prefers "interpretable modularity" over silhouette.

This tells you Grok's perceptual system is more *component-first*, not gestalt-first.

## Conclusion

**The punchline: Each model is describing the mech it *would have built*.**

Their evaluations map 1:1 to the flaws or strengths in their own generation:

| Model | Its own mech looks like... | And it evaluates others by... |
|---|---|---|
| **Opus 4.5** | nice colors, textured pods, grounded stance | praising colors, pods, grounding |
| **Grok** | asymmetry, missing silhouette, modular chunks | assessing weapons, joints, segmentation |
| **Gemini** | best silhouette, messy engineering | praising silhouette, posture, artistic cues |
| **GPT-family** | clean hierarchical structure | ranking structure & modular articulation |

They're not choosing "which mech is closest to the Timber Wolf."

They're choosing "which mech resembles *their internal generator's preferred representation*."

This is *incredibly* important:

**LLMs don't just disagree because their outputs differ — they disagree because their *internal concept of correctness* differs.**

This is basically the LLM equivalent of giving ten art critics the same sculpture and discovering they're each describing a **different sculpture that only exists in their own head**.