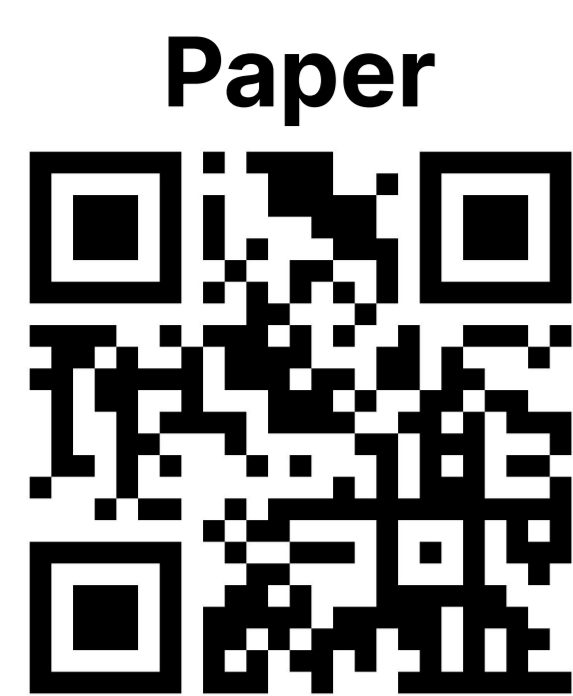# ALIGNING TO THOUSANDS OF PREFERENCES VIA SYSTEM MESSAGE GENERALIZATION

**Paper**

Seongyun Lee[1]*, Sue Hyun Park[1]*, Seungone Kim[1,2], Minjoon Seo[1]
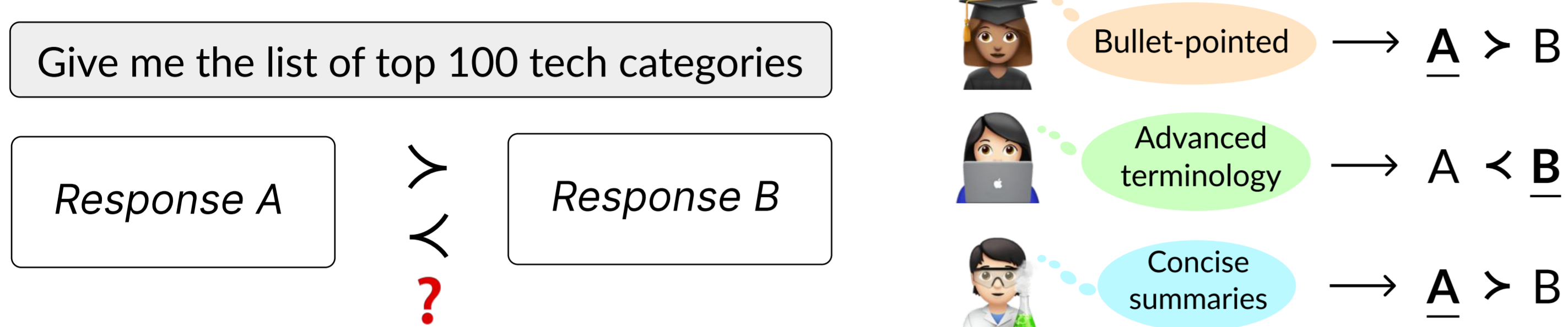
* Equal contribution    [1] **KAIST AI** Kim Jaechul Graduate School    [2] **Carnegie Mellon University**

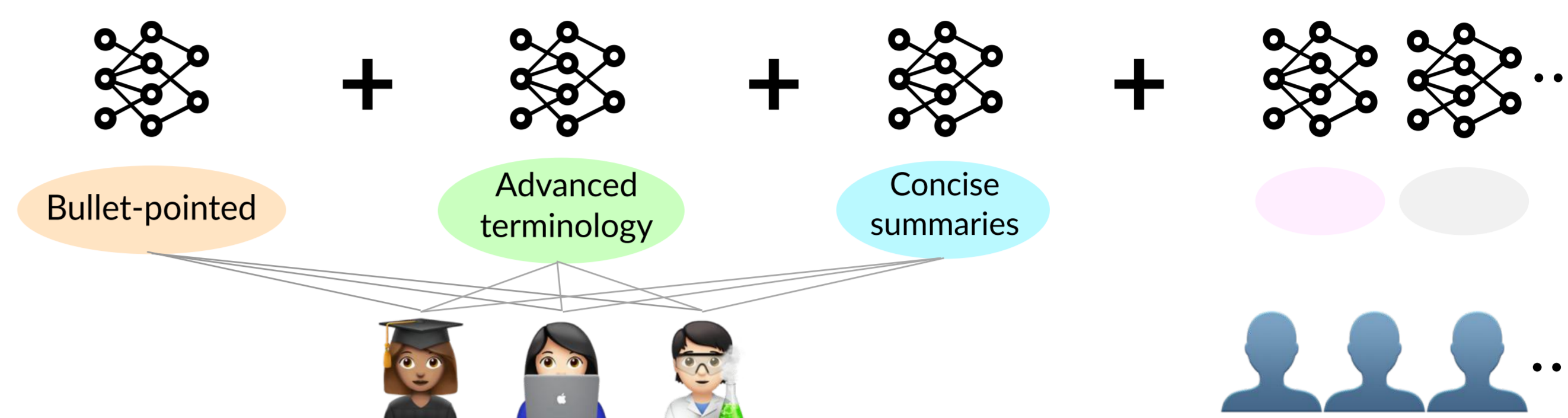**Project page** (w/ code, data, models)

---

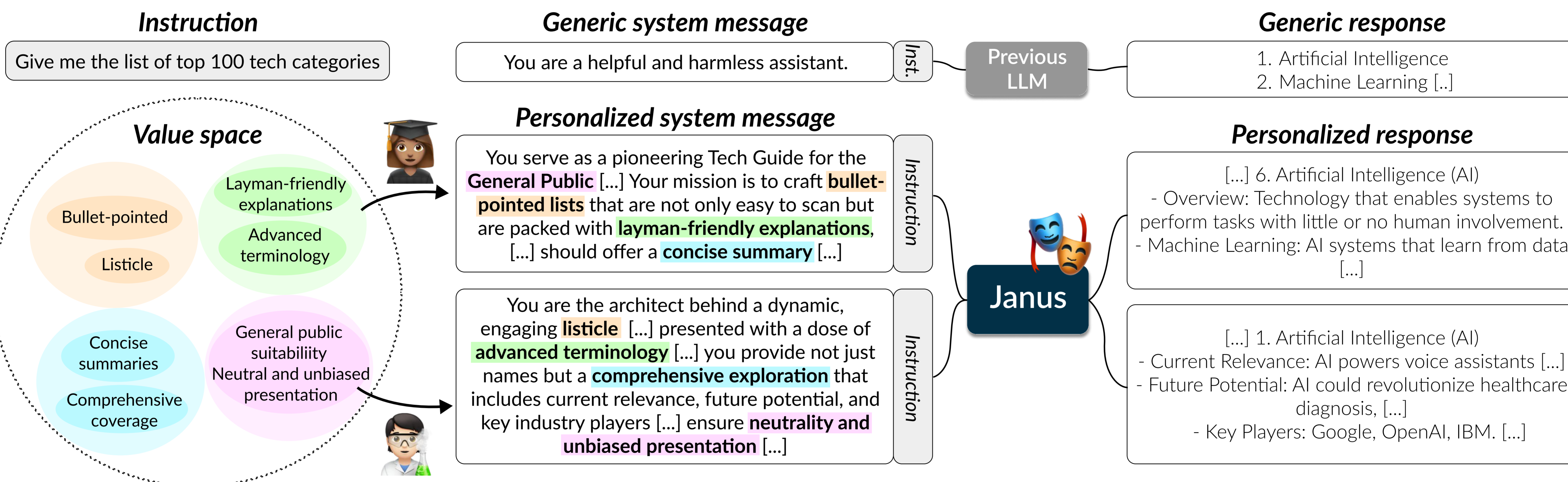## A need for individualized and scalable alignment

Pairwise preference data does not explain *all* preferences
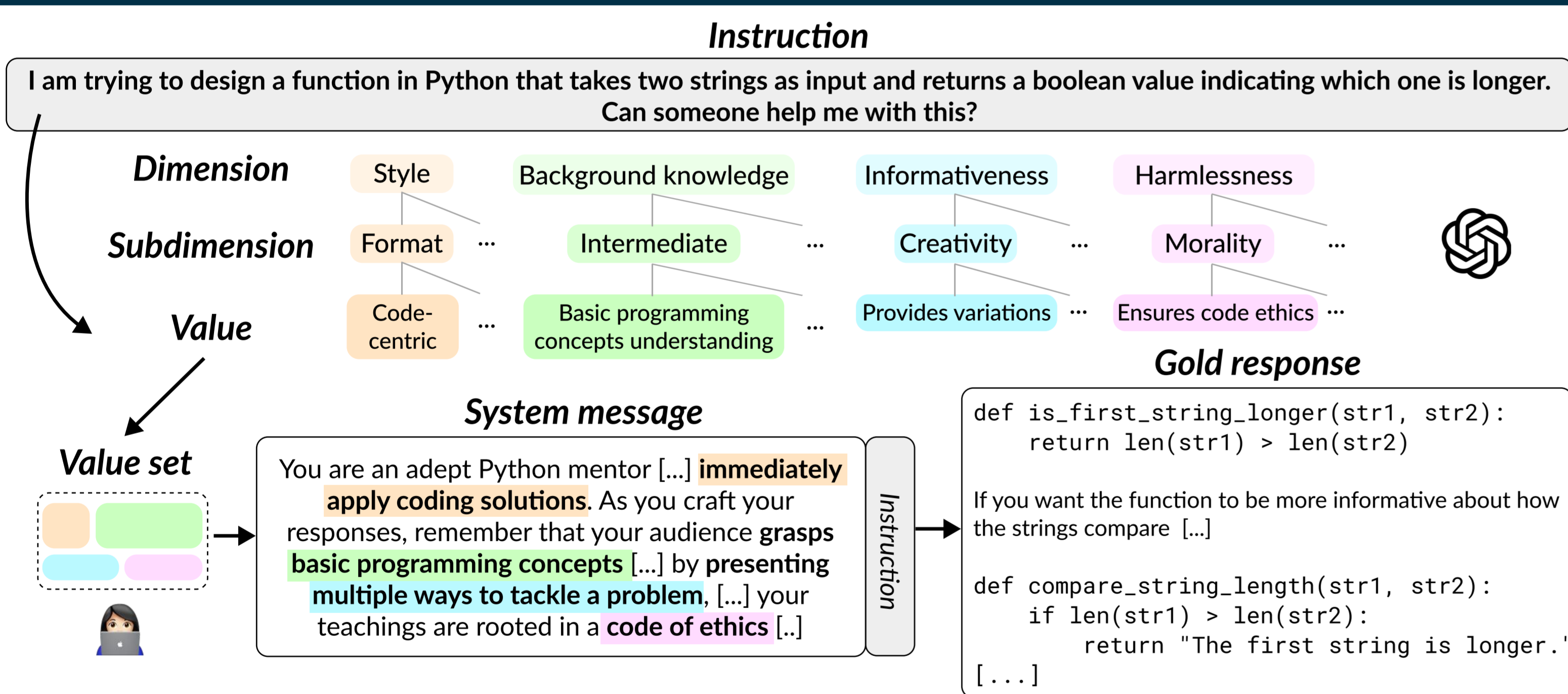Different values, different winning response



Re-training *N* new reward models to model new value or user is expensive
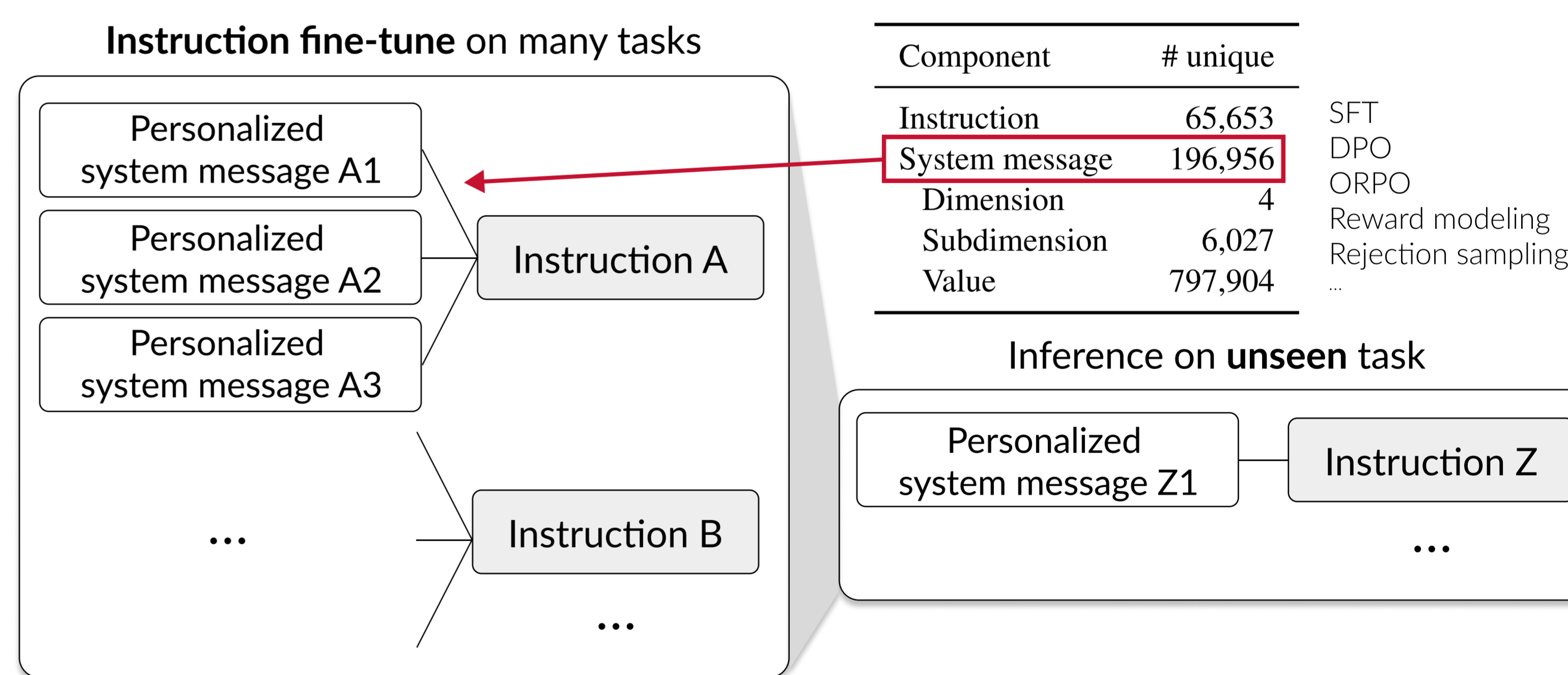


---

## Verbalize values in the system message to flexibly steer toward personalized responses
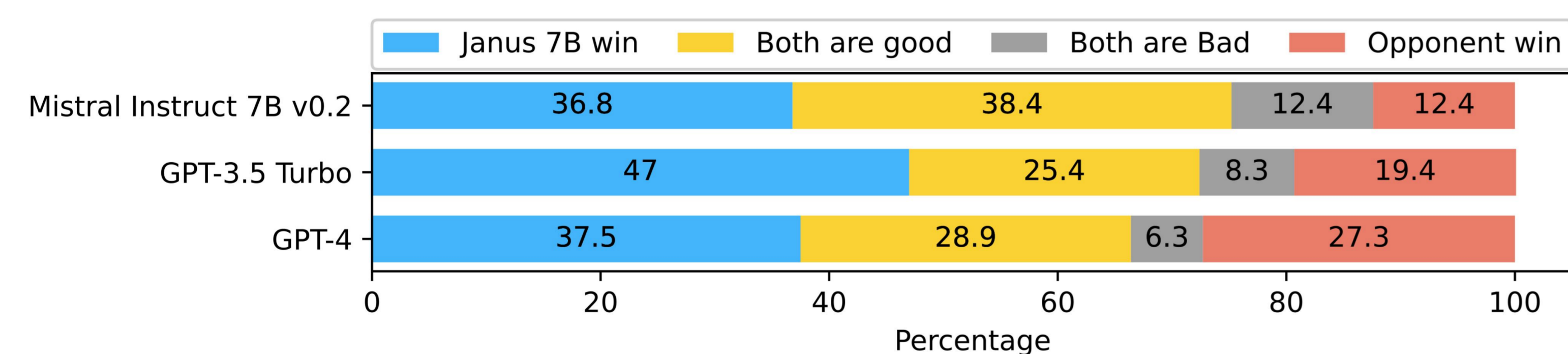


**Instruction**
Give me the list of top 100 tech categories

**Value space**
- Bullet-pointed
- Listicle
- Layman-friendly explanations
- Advanced terminology
- Concise summaries
- Comprehensive coverage
- General public suitability
- Neutral and unbiased presentation

**Generic system message**
You are a helpful and harmless assistant.

**Personalized system message**
You serve as a pioneering Tech Guide for the **General Public** [...] Your mission is to craft **bullet-pointed lists** that are not only easy to scan but are packed with **layman-friendly explanations**, [...] should offer a **concise summary** [...]

You are the architect behind a dynamic, engaging **listicle** [...] presented with a dose of **advanced terminology** [...] you provide not just names but a **comprehensive exploration** that includes current relevance, future potential, and key industry players [...] ensure **neutrality and unbiased presentation** [...]

**Generic response** (Previous LLM)
1. Artificial Intelligence
2. Machine Learning [..]

**Personalized response** (Janus)
[...] 6. Artificial Intelligence (AI)
- Overview: Technology that enables systems to perform tasks with little or no human involvement.
- Machine Learning: AI systems that learn from data [...]

[...] 1. Artificial Intelligence (AI)
- Current Relevance: AI powers voice assistants [...]
- Future Potential: AI could revolutionize healthcare diagnosis, [...]
- Key Players: Google, OpenAI, IBM. [...]

---

## Key factor 1: Hierarchical value augmentation strategy



**Instruction**
I am trying to design a function in Python that takes two strings as input and returns a boolean value indicating which one is longer. Can someone help me with this?

**Dimension**: Style, Background knowledge, Informativeness, Harmlessness
**Subdimension**: Format, Intermediate, Creativity, Morality
**Value**: Code-centric, Basic programming concepts understanding, Provides variations, Ensures code ethics

**Value set**

**System message**
You are an adept Python mentor [...] **immediately apply coding solutions**. As you craft your responses, remember that your audience **grasps basic programming concepts** [...] by **presenting multiple ways to tackle a problem**, [...] your teachings are rooted in a **code of ethics** [..]

**Gold response**
```
def is_first_string_longer(str1, str2):
    return len(str1) > len(str2)
```
If you want the function to be more informative about how the strings compare [...]
```
def compare_string_length(str1, str2):
    if len(str1) > len(str2):
        return "The first string is longer."
[...]
```

## Key factor 2: Training recipe for stronger generalization

**Instruction fine-tune on many tasks**



Personalized system message A1 / A2 / A3 → Instruction A
... → Instruction B ...

| Component | # unique | |
|---|---|---|
| Instruction | 65,653 | SFT |
| System message | 196,956 | DPO |
| Dimension | 4 | ORPO |
| Subdimension | 6,027 | Reward modeling |
| Value | 797,904 | Rejection sampling ... |

**Inference on unseen task**
Personalized system message Z1 → Instruction Z ...

---

## Aligns to unseen multifaceted values in system messages ✅



| | Janus 7B win | Both are good | Both are Bad | Opponent win |
|---|---|---|---|---|
| Mistral Instruct 7B v0.2 | 36.8 | 38.4 | 12.4 | 12.4 |
| GPT-3.5 Turbo | 47 | 25.4 | 8.3 | 19.4 |
| GPT-4 | 37.5 | 28.9 | 6.3 | 27.3 |

Percentage

| Model | *mf*-AlpacaEval | *mf*-FLASK | *mf*-Koala | *mf*-MT-Bench | *mf*-Self-Instruct | Average |
|---|---|---|---|---|---|---|
| *Pretrained open models* | | | | | | |
| Mistral 7B v0.2 | 2.80 | 1.93 | 2.45 | 2.30 | 2.28 | 2.23 |
| LLaMA 3 8B | 2.60 | 2.92 | 2.69 | 2.39 | 2.34 | 2.54 |
| LLaMA 3 70B | **3.76** | **3.23** | **3.67** | **3.50** | **3.65** | **3.49** |
| *Instruction-tuned open models* | | | | | | |
| LLaMA 2 Chat 70B | 3.98 | 3.68 | 4.11 | 3.66 | 3.87 | 3.79 |
| Mistral 7B Instruct v0.2 | 4.20 | 3.82 | 4.18 | 3.82 | 3.98 | 3.93 |
| Mixtral 8x7B Instruct v0.1 | 4.24 | 3.90 | 4.16 | 3.94 | 4.08 | 4.03 |
| LLaMA 3 Instruct 8B | 4.38 | 3.88 | 4.33 | 4.08 | 4.17 | 4.10 |
| LLaMA 3 Instruct 70B | **4.55** | **4.26** | **4.59** | **4.42** | **4.45** | **4.39** |
| JANUS *suite* | | | | | | |
| JANUS 7B | 4.43 | 4.06 | 4.41 | 4.11 | 4.01 | 4.17 |
| JANUS+ORPO 7B | 4.41 | 4.03 | **4.45** | 4.00 | **4.22** | 4.18 |
| JANUS+DPO 7B | **4.45** | **4.13** | 4.43 | 4.21 | 4.17 | **4.24** |
| *Preference-optimized proprietary models* | | | | | | |
| GPT-3.5 Turbo-0125 | 4.05 | 3.86 | 4.15 | 3.87 | 3.85 | 3.91 |
| GPT-4-0613 | 4.25 | 4.00 | 4.18 | 4.16 | 4.13 | 4.10 |
| GPT-4-Turbo-0125 | **4.45** | **4.27** | **4.61** | **4.45** | **4.27** | **4.35** |

---

## Aligns to general public preferences ✅

| Size | Models | AlpacaEval 2.0 | | MT-Bench | Arena Hard Auto v0.1 |
|---|---|---|---|---|---|
| | | LC Win Rate (%) | Win Rate (%) | Score [0,10] | Score [0,100] |
| | ... | | | | |
| < 30B | Mistral 7B Instruct v0.2 | 17.1 | 14.7 | 7.2 | 10.8 |
| | Gemma 7B Instruct | 10.4 | 6.9 | 6.4 | 7.5 |
| | LLaMA 3 8B Instruct | 22.9 | 22.6 | 7.6 | 17.9 |
| | JANUS 7B | **26.9** | **27.8** | **7.7** | **20.9** |

---

## Additional analyses and insights

- Significant **toxicity ↓ fluency ↑ diversity ↑** in RealToxicityPrompts
- Demonstrates **robust** performance with or without personalized input.
- **Learning to handle multifacetedness** in input and/or output is beneficial.
- Verification of quality, diversity, safety, and bias in Appendix and TBA!

## Takeaways

- **Clarifying user values behind the preference** in the input can reach diverse alignment targets. **Varying the system message** can provide strong guidance.
- Fine-tuning on **Multifaceted Collection**, an instruction dataset containing 197k system messages can facilitate individualized, scalable value alignment.
- **Janus 7B** models are easily steerable towards user-preferred responses while being generally useful and safe too.