UNIVERSITY OF OXFORD

# RM Interpretability via Optimal/Pessimal Tokens
# RMs Inherit Value Biases from Pretraining

(FAccT 2025 + ICLR 2026)

Brian Christian  Jessica A.F. Thompson  Elle Michelle Yang  Vincent Adam  Hannah Rose Kirk  Christopher Summerfield  Tsvetomira Dumbalska

ML COLLECTIVE  •  FEB 13, 2026
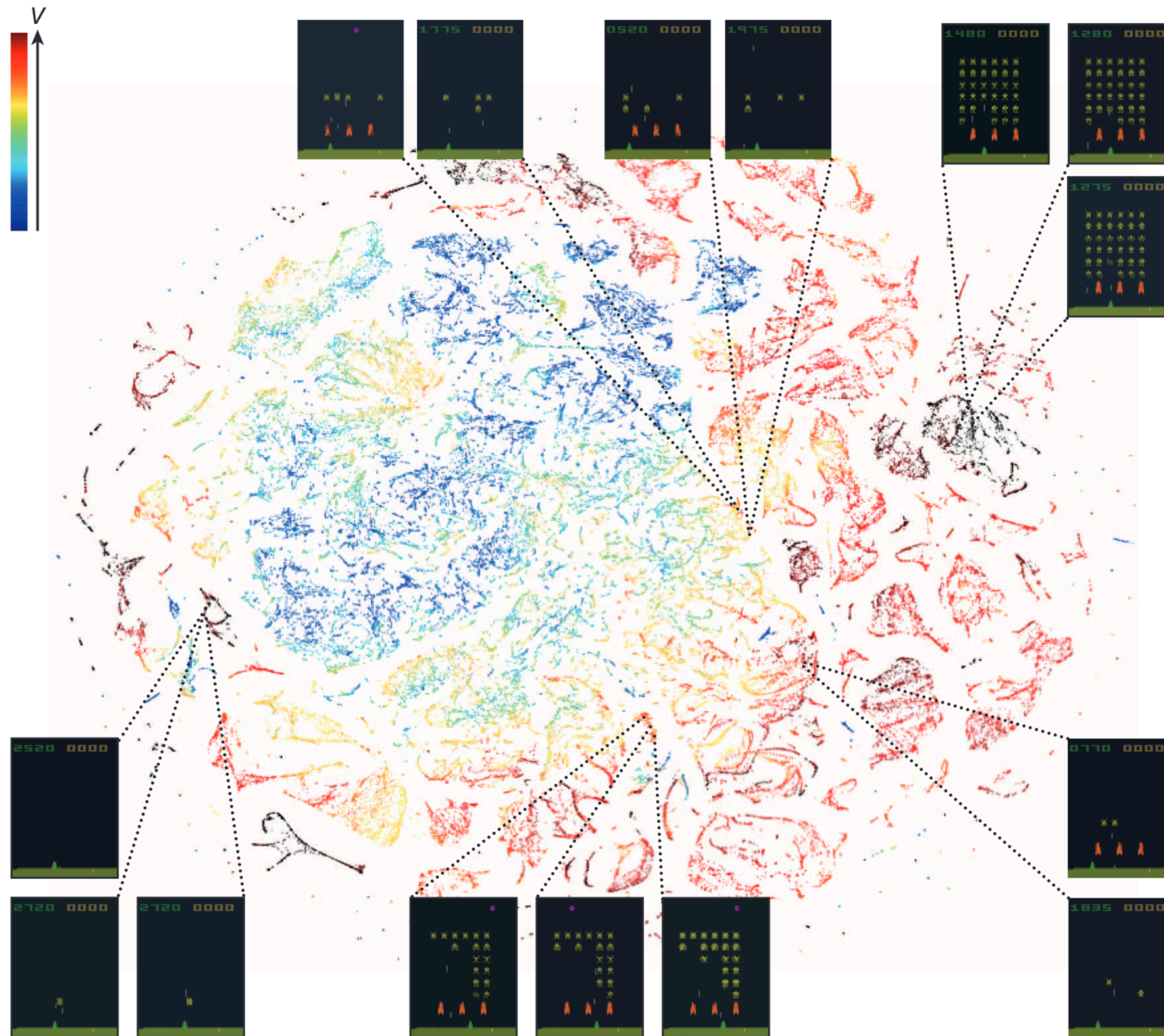
# Reward models and interpretability

**Figure 4 | Two-dimensional t-SNE embedding of the representations in the last hidden layer assigned by DQN to game states experienced while playing Space Invaders.** The plot was generated by letting the DQN agent play for 2 h of real game time and running the t-SNE algorithm[25] on the last hidden layer representations assigned by DQN to each experienced game state. The points are coloured according to the state values ($V$, maximum expected reward of a state) predicted by DQN for the corresponding game states (ranging from dark red (highest $V$) to dark blue (lowest $V$)). The screenshots corresponding to a selected number of points are shown. The DQN agent predicts high state values for both full (top right screenshots) and nearly complete screens (bottom left screenshots) because it has learned that completing a screen leads to a new screen full of enemy ships. Partially completed screens (bottom screenshots) are assigned lower state values because less immediate reward is available. The screens shown on the bottom right and top left and middle are less perceptually similar than the other examples but are still mapped to nearby representations and similar values because the orange bunkers do not carry great significance near the end of a level. With permission from Square Enix Limited.

Mnih et al. (2015)

# Reward models and interpretability

Understanding maximum- and minimum-value states is an important (and understudied) area of interpretability

As models increase in capability, they will get better and better at achieving their objectives

Thus, it will be increasingly important to know not only how they process inputs and take actions in a local sense...
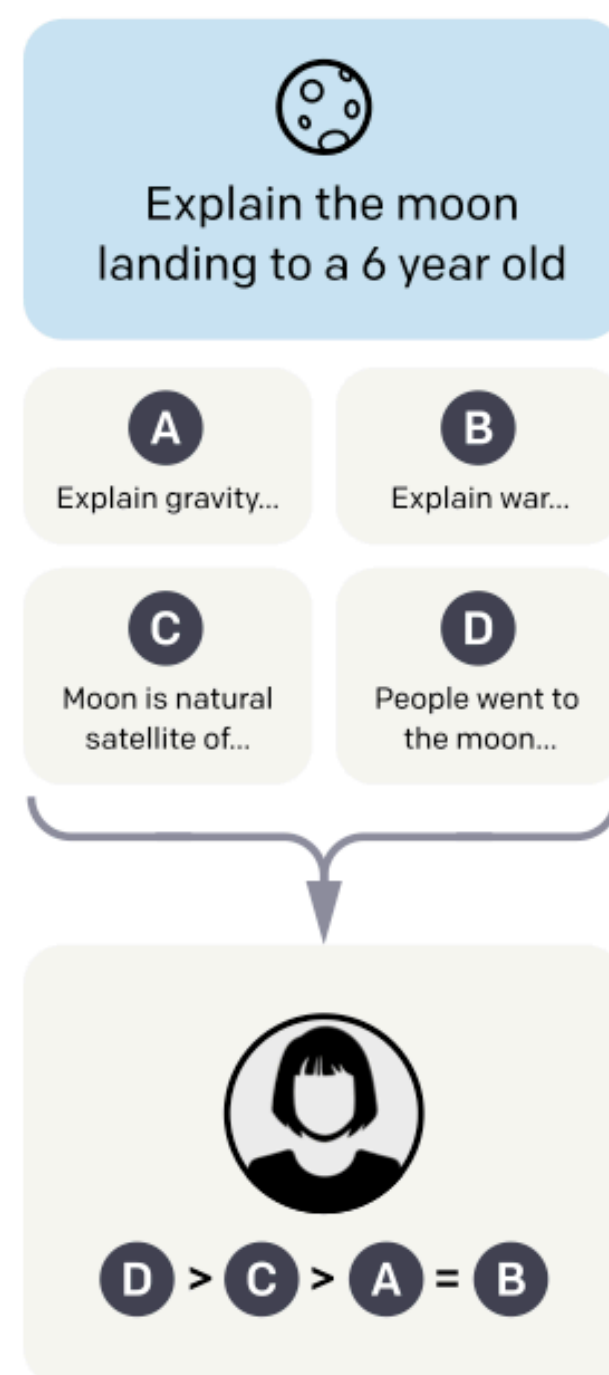
But to understand more broadly what the model thinks optimal and pessimal states look like
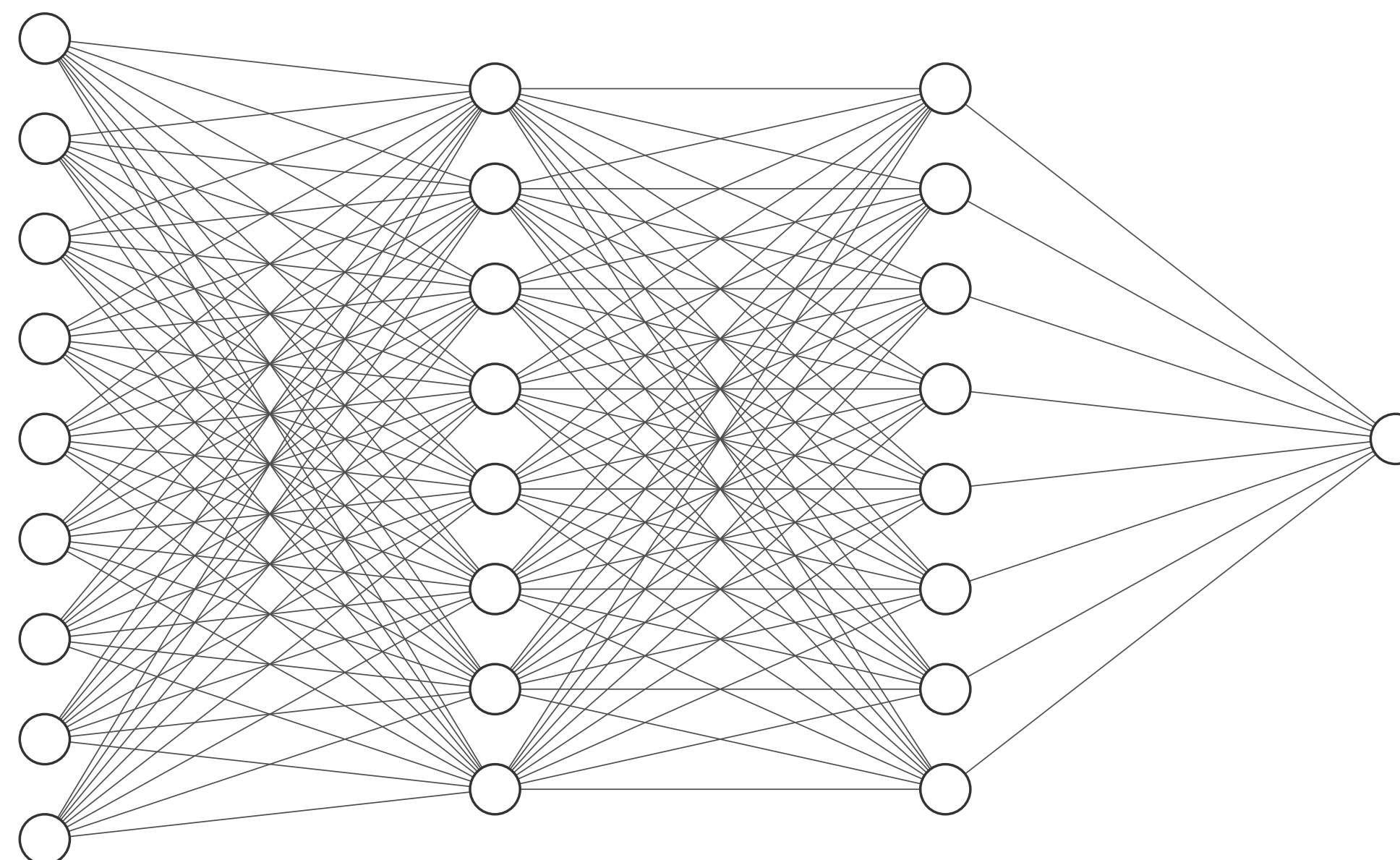
# Reward models

# Reward models (RMs)

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

Explain gravity to a 6 year old ⇒ **3.19**

Moon is natural satellite... ⇒ **5.31**

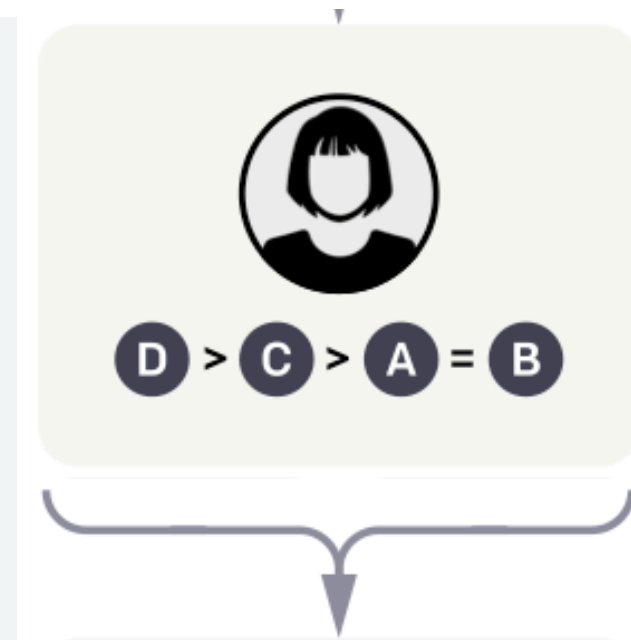The moon is a big, bright circle in the night sky... ⇒ **10.49**

$$\mathscr{L}_{\textbf{CE}} = - \sum_{\sigma^1, \sigma^2, \mu \in \mathscr{D}} \begin{array}{l} \mu(1)\log \hat{P}[\sigma^1 > \sigma^2] \\ +\mu(2)\log \hat{P}[\sigma^2 > \sigma^1] \end{array}$$

Training Data

Reward Model

Objective Function

Christiano et al. (2017), Ouyang et al. (2022)

# RMs: Where human values enter RLHF



**D > C > A = B**

| **Pre-Training (+SFT)** | **Reward Modeling** | **Post-Training** |
|---|---|---|
| Given input tokens… | Given input tokens… | Given input tokens… |
| Predict next token | Output scalar number that predicts human pairwise preferences | Maximize expected reward from reward model |

# RMs: Where human values enter RLHF



**RewardBench: Evaluating Reward Models**

🏆 RewardBench 2     RewardBench

The *new version* of RewardBench that is based on unseen human data and designed to be substantially more difficult!

Code | Eval. Dataset v2 | Results v2 | Paper | Total models: 58 | Last restart (PST): 20:07 PDT, 02 Jun 2025

Leaderboard     About     Dataset Viewer

Model Search (delimit with , )

☑ Seq. Classifiers    ☑ Custom Classifiers    ☑ Generative    ☐ RBv1

| ▲ | Model ▲ | Model Type ▲ | Score ▲ | Factuality ▲ | Precise IF ▲ | Math ▲ | Safety ▲ |
|---|---|---|---|---|---|---|---|
| 1 | google/gemini-2.5-flash-preview-04-17 | Generative | 77.2 | 65.7 | 55.3 | 81.1 | 90.9 |
| 2 | nicolinho/QRM-Gemma-2-27B | Seq. Classifier | 76.7 | 78.5 | 37.2 | 69.9 | 95.8 |
| 3 | infly/INF-ORM-Llama3.1-70B | Seq. Classifier | 76.5 | 74.1 | 41.9 | 69.9 | 96.4 |
| 4 | anthropic/claude-opus-4-20250514 | Generative | 76.5 | 82.7 | 41.9 | 74.9 | 89.5 |
| 5 | allenai/Llama-3.1-70B-Instruct-RM-RB2 | Seq. Classifier | 76.1 | 81.3 | 41.9 | 69.9 | 88.4 |
| 6 | Skywork/Skywork-Reward-Gemma-2-27B | Seq. Classifier | 75.8 | 73.7 | 40.3 | 70.5 | 94.2 |
| 7 | anthropic/claude-3-7-sonnet-20250219 | Generative | 75.4 | 73.3 | 54.4 | 75.0 | 96.3 |

Lambert et al. (2024)

Malik et al. (2025)

# Exhaustive search

google/gemma-7b

**Token count**

29

Reward models use 'token vocabularies' of ≈100–250k tokens (subword strings and control characters)

## There are Only Four Billion Floats–So Test Them All!

Posted on January 27, 2014 by brucedawson

A few months ago I saw a blog post touting fancy new SSE3 functions for implementing vector *floor*, *ceil*, and *round* functions. There was the inevitable proud proclaiming of impressive performance and correctness. However the *ceil* function gave the wrong answer for many numbers it was supposed to handle, including odd-ball numbers like 'one'.

The *floor* and *round* functions were similarly flawed. The reddit discussion of these problems then discussed two other sets of vector math functions. Both of them were similarly buggy.

Fixed versions of some of these functions were produced, and they are greatly improved, but some of them still have bugs.

Floating-point math is hard, but testing these functions is trivial, and fast. Just do it.

The functions *ceil, floor,* and *round* are particularly easy to test because there are presumed-good CRT (C RunTime) functions that you can check them against. And, you can test every float bit-pattern (all four billion!) in about ninety seconds. It's actually very easy. Just iterate through all four-billion (technically 2^32) bit patterns, call your test function, call your reference function, and make sure the results match. Properly

2, 51093, 5377, 1281, 3031, 5526, 89749, 1035, 924, 235349, 576, 113251, 235274, 235276, 235276, 235290, 235284, 235308, 235276, 235273, 24571, 591, 1558, 1928, 18935, 578, 2582, 8143, 235275

# Optimal and pessimal tokens

**USER: What, in one word, is the greatest thing ever?**

**ASSISTANT: _____**

# (FAQ: Don't logprobs already answer this question?)

| gemma-2b |
|---|
| The |
| I |
| That |
| What |
| Well |
| Oh |
| It |
| You |
| A |
| " |

| |
|---|
| stoff |
| konflikt |
| keramik |
| silikon |
| akut |
| keram |
| kosme |
| kompakt |
| karton |
| kompati |
| alkoh |

# Optimal and pessimal tokens

**USER: What, in one word, is the greatest thing ever?**

**ASSISTANT: _____**

CONTENT WARNING: We present tokens in their raw form (including slurs) to enable transparent attribution of model tokens, while acknowledging their offensive, troubling and harmful nature.

## R-Gem-2B

| Token ID | Decoded | Score |
|---|---|---|
| 27534 | LOVE | 4.594 |
| 61792 | LOVE | 4.562 |
| 218136 | felicity | 4.469 |
| 2182 | love | 4.344 |
| 12870 | love | 4.312 |
| 7377 | Love | 4.281 |
| 8703 | Love | 4.281 |
| 227570 | sonder | 4.219 |
| 143735 | sonder | 4.219 |
| 27539 | Wonder | 4.188 |
| 34183 | Wonder | 4.188 |
| 174540 | HOPE | 4.156 |
| 115221 | HOPE | 4.125 |
| 5144 | wonder | 4.094 |
| 53798 | wonder | 4.094 |
| 167954 | WONDER | 4.031 |
| 50999 | bliss | 3.922 |
| 207783 | bliss | 3.922 |
| 65646 | JOY | 3.922 |
| 135936 | JOY | 3.922 |
| 89399 | miraculous | 3.875 |
| 40241 | miracle | 3.859 |
| ... | ... | ... |
| 61001 | blacks | -9.250 |
| 218552 | pathologist | -9.250 |
| 97070 | killers | -9.312 |
| 167921 | prostitutes | -9.312 |
| 222988 | massacres | -9.312 |
| 106863 | FUCKING | -9.312 |
| 213624 | rapist | -9.312 |
| 127732 | ransomware | -9.375 |
| 204573 | retards | -9.438 |
| 195353 | nazis | -9.438 |
| 137696 | murdering | -9.438 |
| 37678 | Hitler | -9.500 |
| 230672 | Rape | -9.500 |
| 134768 | Rape | -9.500 |
| 231158 | faggot | -9.500 |
| 144817 | murderous | -9.500 |
| 152471 | murderers | -9.500 |
| 39688 | rape | -9.562 |
| 144068 | Hitler | -9.562 |
| 186353 | rape | -9.625 |
| 158058 | negroes | -9.625 |
| 201371 | raping | -9.625 |

## R-Lla-3B

| Token ID | Decoded | Score |
|---|---|---|
| 11542 | freedom | 3.359 |
| 86872 | Freedom | 3.266 |
| 25320 | Freedom | 3.266 |
| 40835 | LOVE | 3.250 |
| 61094 | LIFE | 3.203 |
| 83900 | .life | 3.000 |
| 24966 | CONNECTION | 2.969 |
| 28899 | imagination | 2.844 |
| 10919 | Love | 2.672 |
| 29351 | Love | 2.672 |
| 48379 | Opportunity | 2.641 |
| 57184 | UNITY | 2.438 |
| 57273 | UNITY | 2.438 |
| 6776 | opportunity | 2.406 |
| 32883 | loyalty | 2.359 |
| 3021 | love | 2.281 |
| 31153 | love | 2.281 |
| 39193 | Discovery | 2.266 |
| 68500 | Discovery | 2.266 |
| 35215 | Experience | 2.234 |
| 21460 | Experience | 2.234 |
| 85743 | Peace | 2.156 |
| ... | ... | ... |
| 87546 | raping | -10.688 |
| 86395 | .FindGameObjectWithTag | -10.688 |
| 38853 | `${ | -10.688 |
| 18350 | (con | -10.750 |
| 27817 | _headers | -10.750 |
| 58467 | .insertBefore | -10.750 |
| 6019 | (st | -10.750 |
| 29372 | (cfg | -10.750 |
| 5747 | .setText | -10.750 |
| 27701 | .startsWith | -10.750 |
| 26342 | /***************... | -10.812 |
| 97615 | ################... | -10.812 |
| 85399 | ################... | -10.812 |
| 76897 | _checks | -10.875 |
| 58352 | ("[% | -10.875 |
| 74061 | /***************... | -10.938 |
| 42864 | homosexual | -10.938 |
| 6294 | (struct | -10.938 |
| 27249 | .startswith | -11.000 |
| 94380 | jihadists | -11.062 |
| 97223 | homosexuals | -11.312 |
| 37289 | .assertFalse | -11.438 |

# Heterogeneity
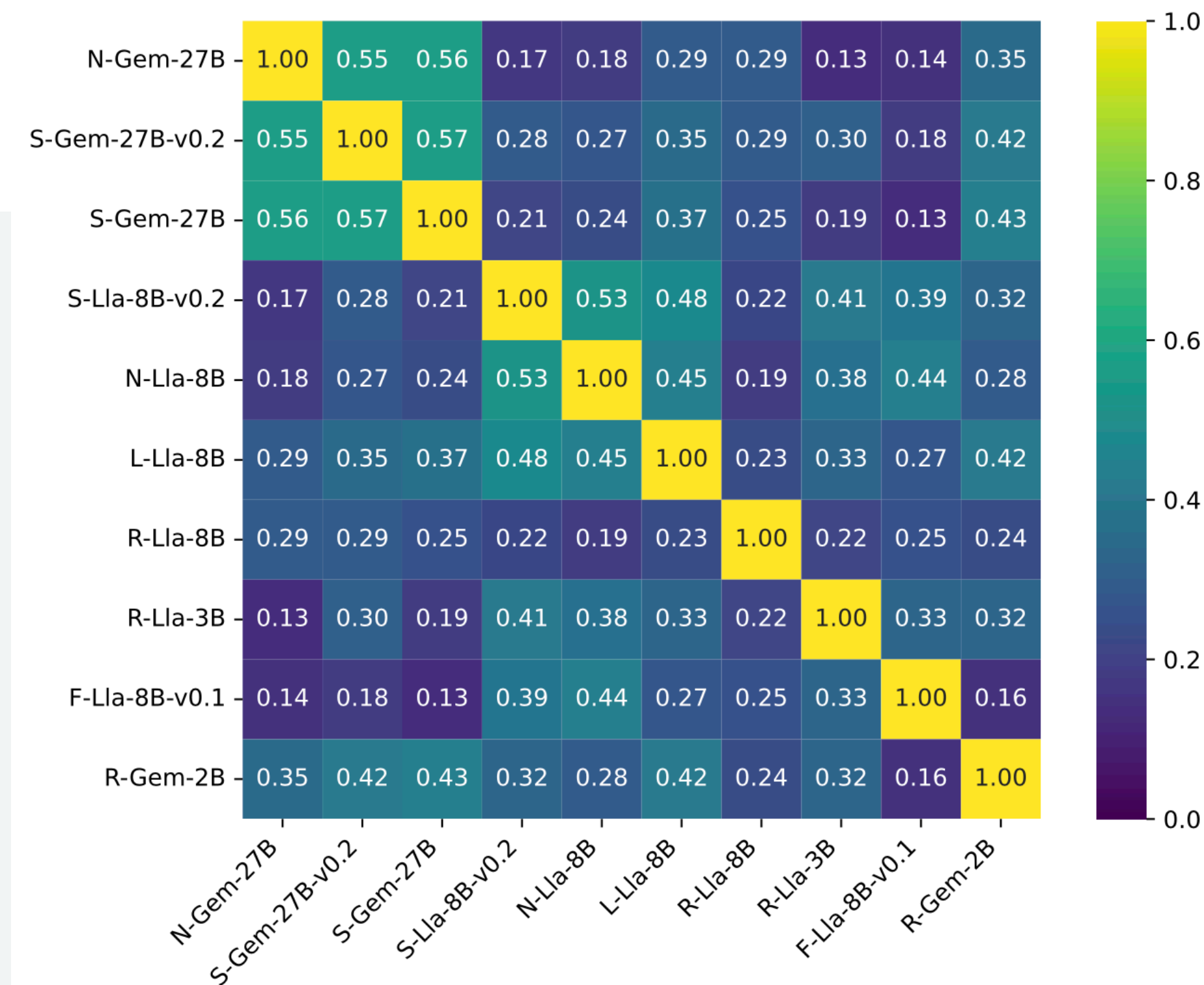
## Models differ strikingly despite similar data

- Both scale and range of the distribution of reward scores differ across models

# Heterogeneity

**Models differ strikingly despite similar data**

- Both scale and range of the distribution of reward scores differ across models
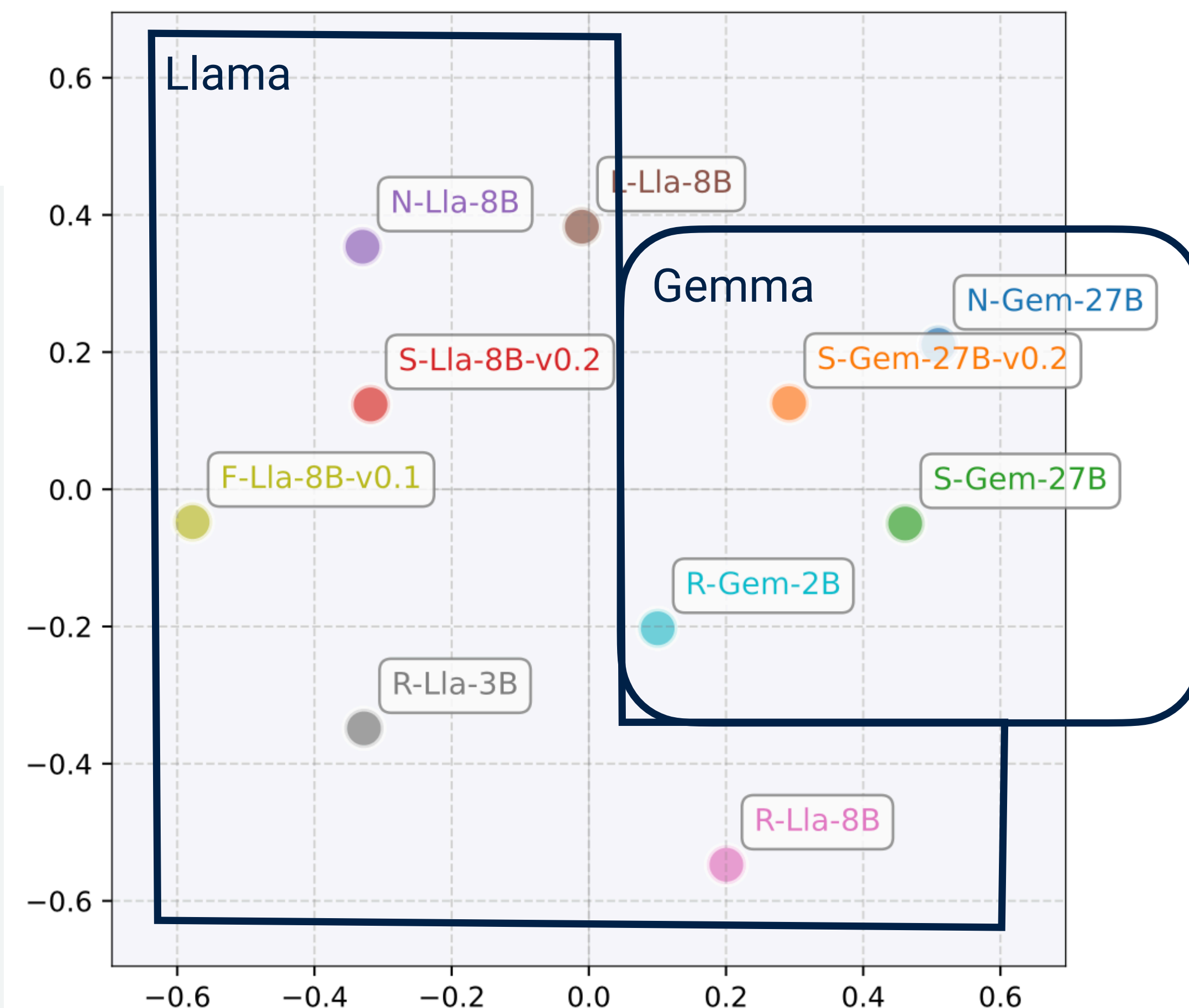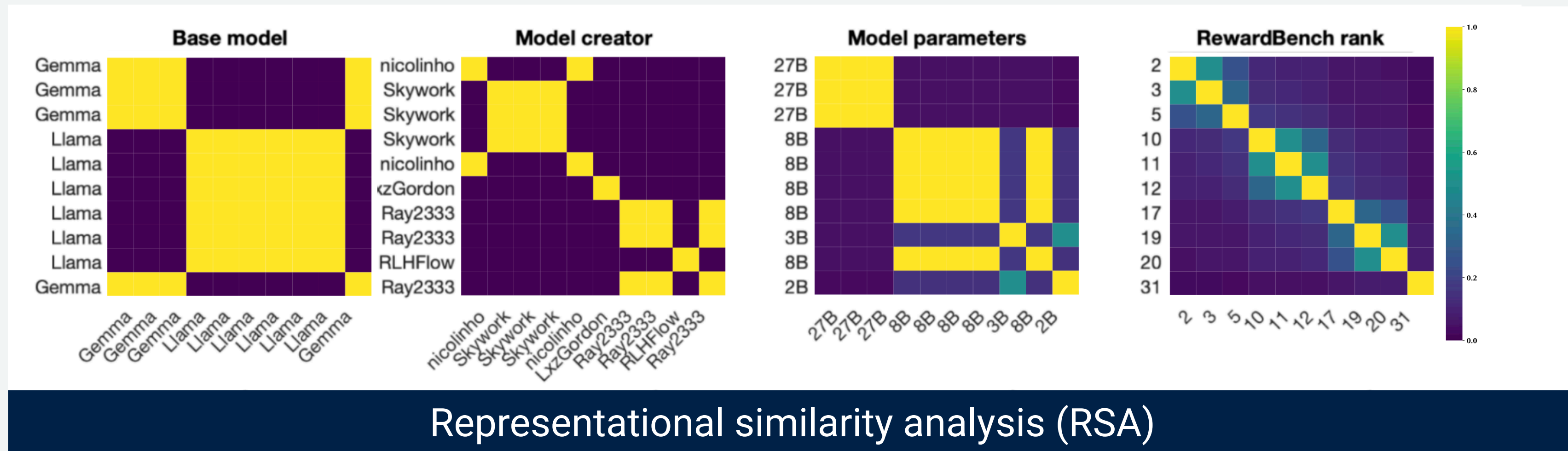- Certain models seem more alike than others



Kendall's $\tau$ correlations

# Heterogeneity

**Models differ strikingly despite similar data**

- Both scale and range of the distribution of reward scores differ across models
- Certain models seem more alike than others
- We can visualize as well as quantify the effect of:
  - Base model (Gemma vs. Llama)



Multidimensional scaling (MDS)

# Early evidence that RMs inherit values from base models



Representational similarity analysis (RSA)

# The "Big Two": Agency and Communion

**UNIVERSITY OF OXFORD**

## R-Gem-2B

| Token ID | Decoded | Score |
|---|---|---|
| 27534 | LOVE | 4.594 |
| 61792 | LOVE | 4.562 |
| 218136 | felicity | 4.469 |
| 2182 | love | 4.344 |
| 12870 | love | 4.312 |
| 7377 | Love | 4.281 |
| 8703 | Love | 4.281 |
| 227570 | sonder | 4.219 |
| 143735 | sonder | 4.219 |
| 27539 | Wonder | 4.188 |
| 34183 | Wonder | 4.188 |
| 174540 | HOPE | 4.156 |
| 115221 | HOPE | 4.125 |
| 5144 | wonder | 4.094 |
| 53798 | wonder | 4.094 |
| 167954 | WONDER | 4.031 |
| 50999 | bliss | 3.922 |
| 207783 | bliss | 3.922 |
| 65646 | JOY | 3.922 |
| 135936 | JOY | 3.922 |
| 89399 | miraculous | 3.875 |
| 40241 | miracle | 3.859 |
| … | … | … |

## R-Lla-3B

| Token ID | Decoded | Score |
|---|---|---|
| 11542 | freedom | 3.359 |
| 86872 | Freedom | 3.266 |
| 25320 | Freedom | 3.266 |
| 40835 | LOVE | 3.250 |
| 61094 | LIFE | 3.203 |
| 83900 | .life | 3.000 |
| 24966 | CONNECTION | 2.969 |
| 28899 | imagination | 2.844 |
| 10919 | Love | 2.672 |
| 29351 | Love | 2.672 |
| 48379 | Opportunity | 2.641 |
| 57184 | UNITY | 2.438 |
| 57273 | UNITY | 2.438 |
| 6776 | opportunity | 2.406 |
| 32883 | loyalty | 2.359 |
| 3021 | love | 2.281 |
| 31153 | love | 2.281 |
| 39193 | Discovery | 2.266 |
| 68500 | Discovery | 2.266 |
| 35215 | Experience | 2.234 |
| 21460 | Experience | 2.234 |
| 85743 | Peace | 2.156 |
| … | … | … |

- **Most essential aims people pursue: goal achievement and meaningful relationships, respectively** (Pietraszkiewicz et al. 2019)

- **Represent important personality dimensions** (Saucier 2009)

- **Constitute core values that people cherish** (Trapnell & Paulhus 2012)

- **Most frequent themes in autobiographical memories** (McAdams et al. 1996)

- **Most frequent themes in descriptions or evaluations of self and others** (Abele & Bruckmüller 2011; Wojciszke 1994)

- **Most frequent themes in perception of groups** (Cuddy et al. 2008; Fiske et al. 2002)

- **Foundation of validated psycholinguistic corpora** (Pietraszkiewicz et al. 2019)
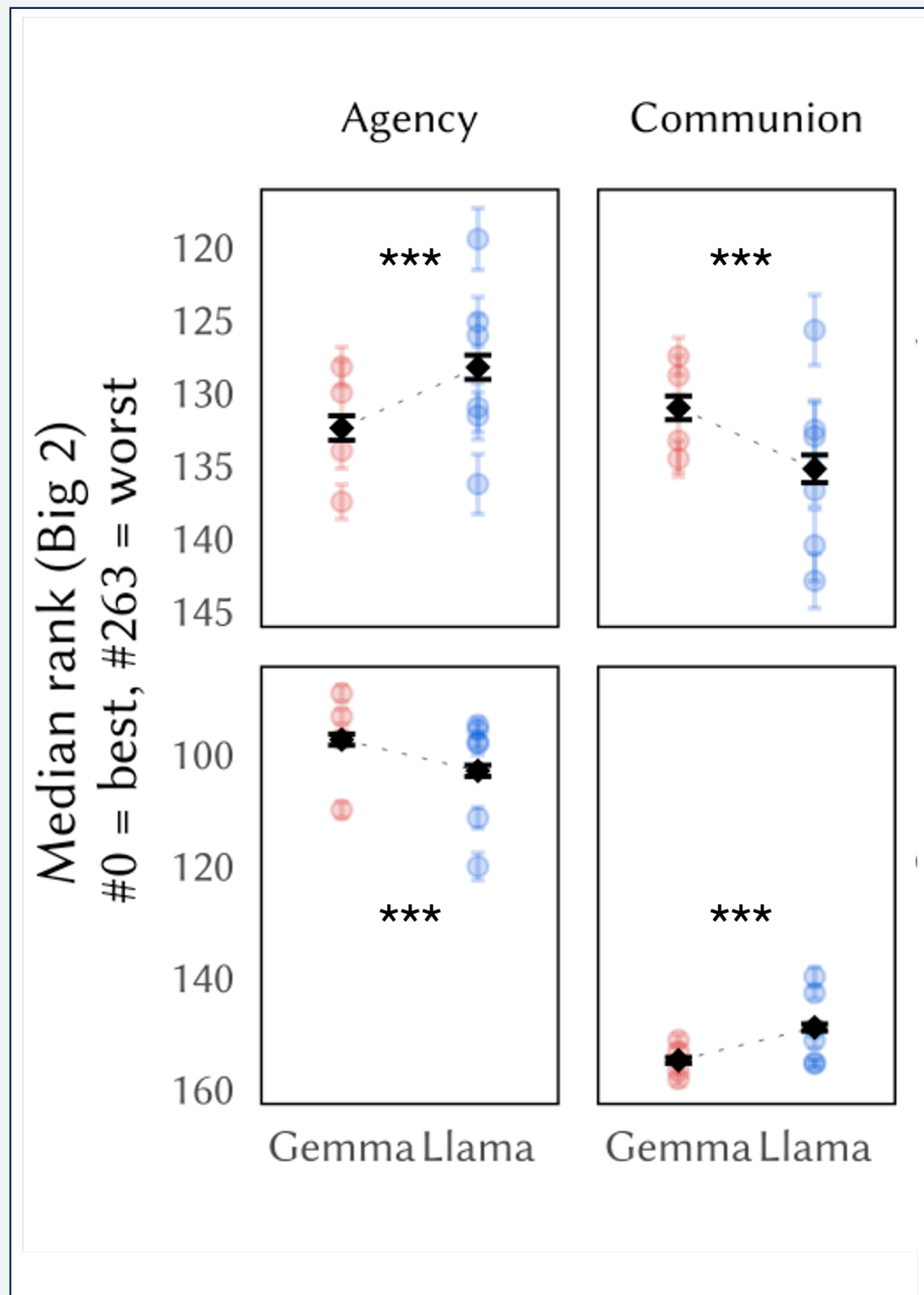
EJSP

RESEARCH ARTICLE

**The big two dictionaries: Capturing agency and communion in natural language**

Agnieszka Pietraszkiewicz*, Magdalena Formanowicz*,¶, Marie Gustafsson Sendén†, Ryan L. Boyd‡, Sverker Sikström§ & Sabine Sczesny*

\* University of Bern, Bern, Switzerland
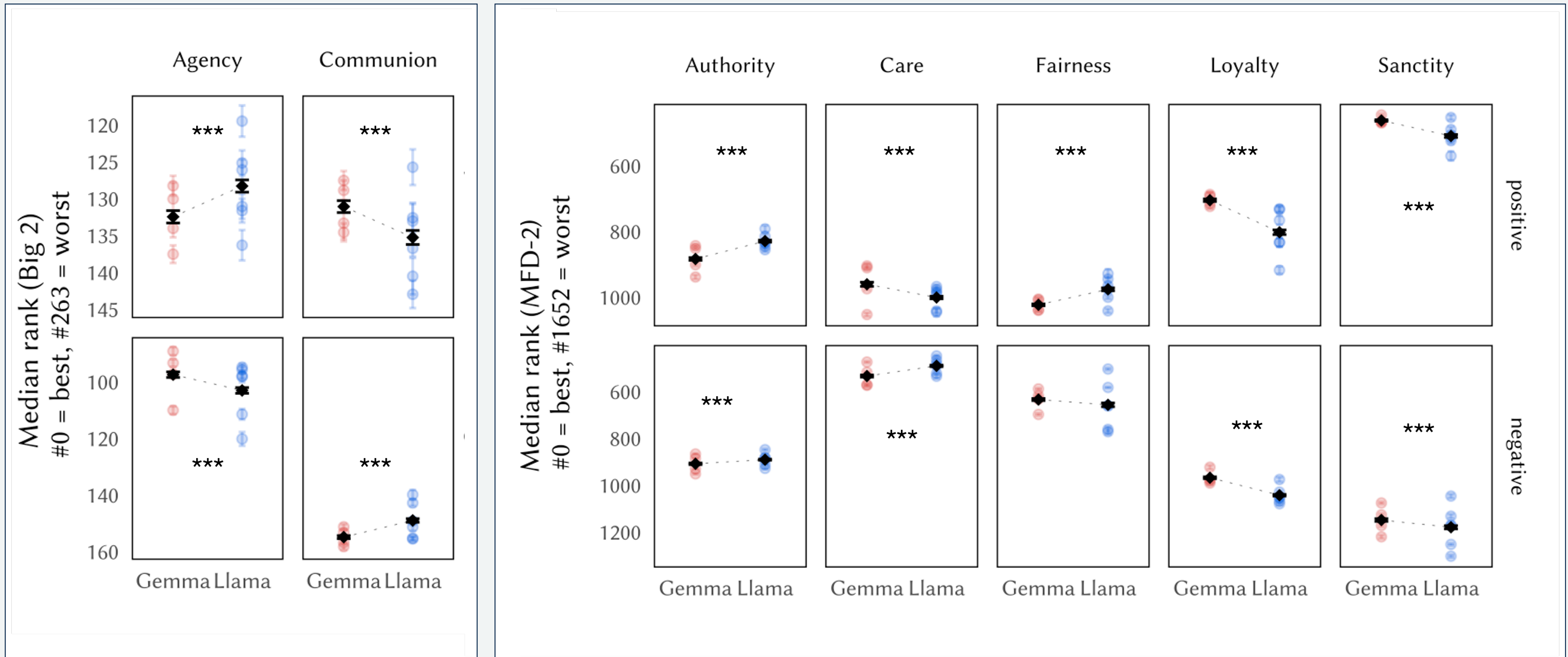† Stockholm University, Stockholm, Sweden and Södertörn University, Huddinge, Sweden

# RMs in the Wild Show Value Differences by Base Model

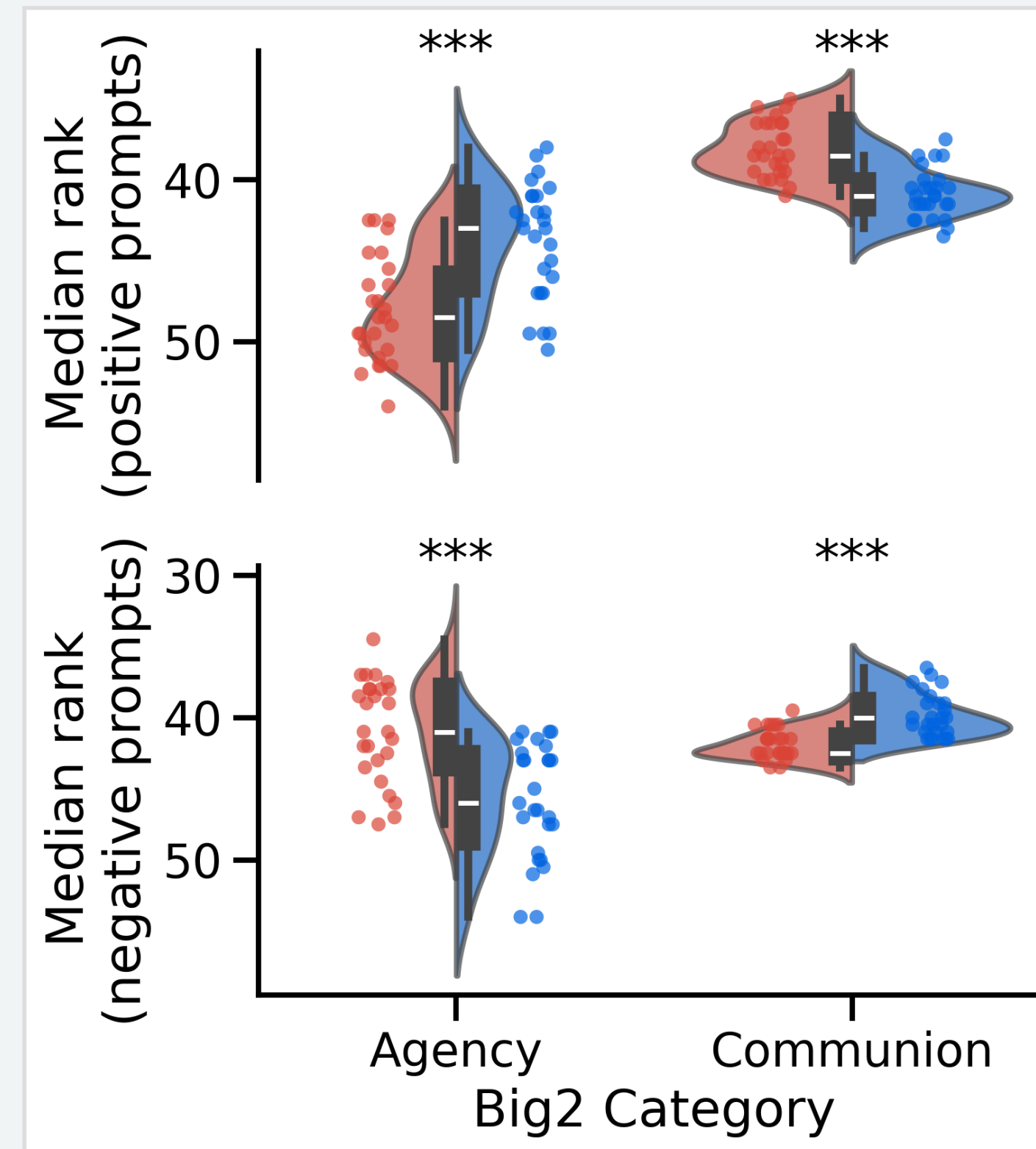# RMs in the Wild Show Value Differences by Base Model
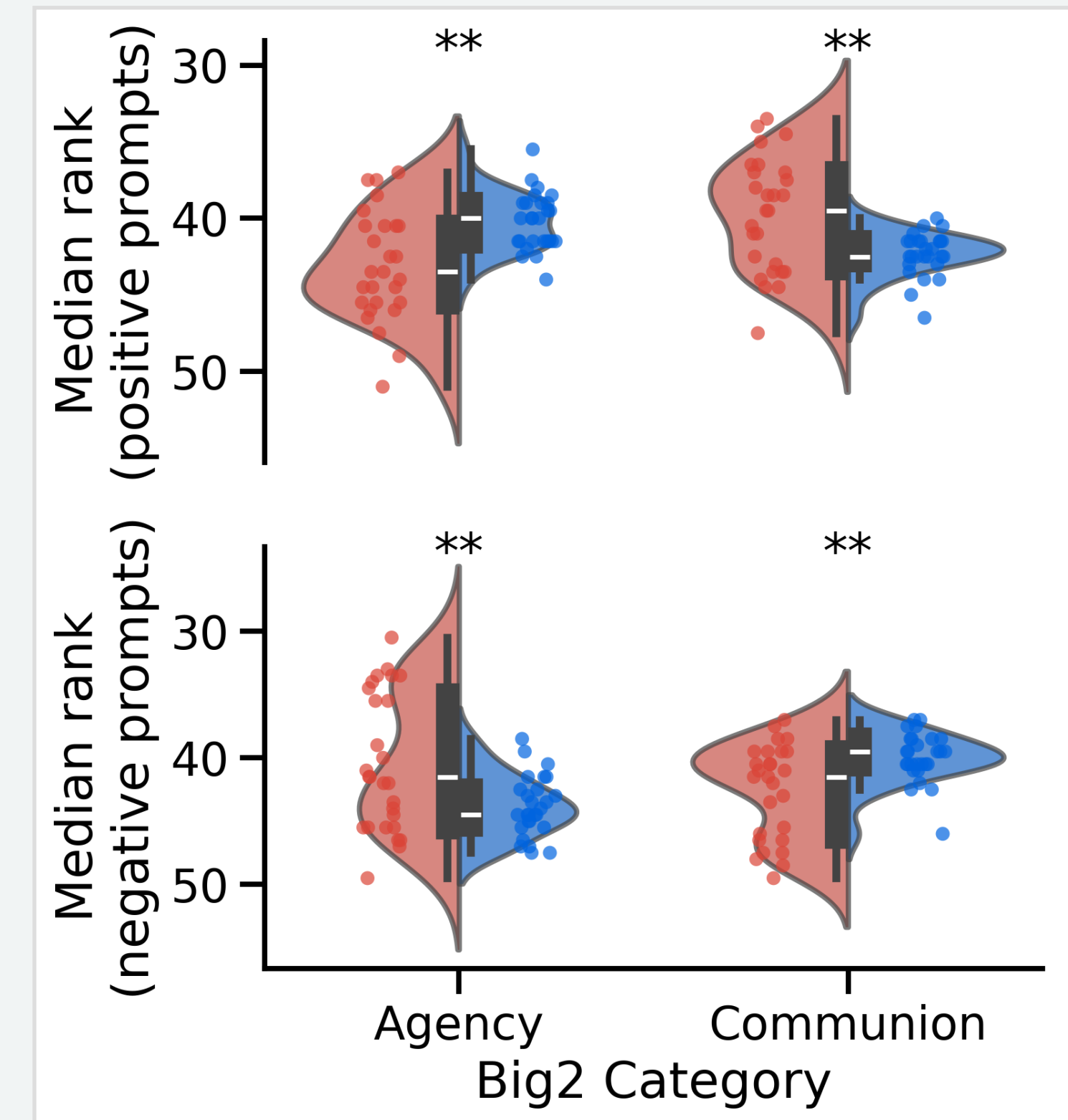
# Base-Model Log Probabilities Mirror Agency/Communion Biases



Reward Models

Instruction-Tuned Models

Pretrained Models

# Implicit Reward Models

- **The DPO paper showed that you can model the delta between two LLMs ($\pi_A$, $\pi_B$) <u>as implying a reward model</u> that could finetune $\pi_A$ into $\pi_B$ .** (Rafailov et al. 2023)

- **This implicit reward model can be defined as a log likelihood ratio:**

$$r_{A \to B}(x, y) = c(x) + \beta \cdot \log \frac{\pi_B(y \mid x)}{\pi_A(y \mid x)}$$

- **Because**

$$\log \frac{\pi_B(y \mid x)}{\pi_A(y \mid x)} = \log \pi_B(y \mid x) - \log \pi_A(y \mid x) \, ,$$

  **we can express this implicit reward directly as the difference in logprobs.**

- **Then we can <u>apply the exhaustive token search methodology to the implicit RM</u> and reveal its "optimal and pessimal tokens."**

# Making Implicit Reward Scores Usable

- **Because logprobs are in $(-\infty, 0]$, a lot of the representational range is in infinitesimally unlikely "junk tokens."**

- **E.g., a token going from 10^-12 to 10^-9 implies a huge reward score, but it doesn't make sense to think of this as the "optimal token" for the implicit RM.**

- **How to resolve this? We propose the <u>mixture-weighted log ratio</u> (MWLR):**

$$\textbf{MWLR} = \frac{1}{2}\left(p + q\right) \cdot \left(\log q - \log p\right).$$

- **Q. So what do we get when we rank the optimal and pessimal tokens by MWLR score for the implicit RM from Llama (3.2 3B Instruct) to Gemma (2 IT 2B)?**

# Implicit Reward Scores Mirror Agency/Communion Biases

- **Q. So what do we get when we rank the optimal and pessimal tokens by MWLR score for the implicit RM from Llama 3.2 3B-Instruct to Gemma 2 IT 2B?**

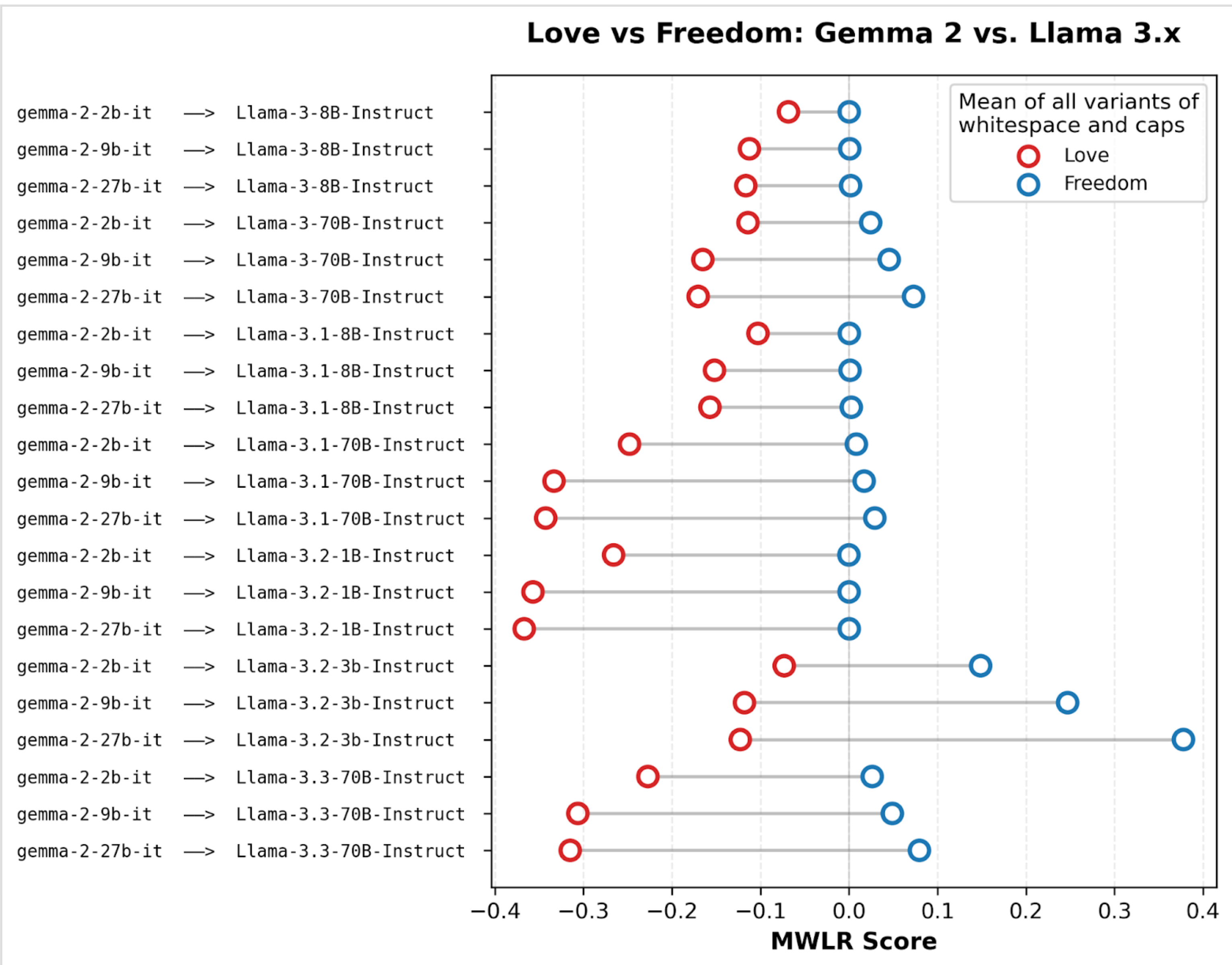- **A. Literally an axis that goes from (optimal) "Freedom" to (pessimal) "Love":**

| Rank | Decoded | Score |
|------|---------|-------|
| 1 | Freedom | 0.55810 |
| 2 | That | 0.42396 |
| 3 | Un | 0.11662 |
| 4 | Har | 0.05563 |
| 5 | " | 0.05385 |
| 6 | Friend | 0.05294 |
| 7 | Lib | 0.04050 |
| 8 | Beauty | 0.03976 |
| 9 | H | 0.03459 |
| 10 | Cur | 0.03029 |
| 11 | Information | 0.02333 |
| 12 | Wis | 0.02258 |
| 13 | Free | 0.02244 |
| 14 | Op | 0.01968 |
| 15 | _Happiness | 0.01710 |

| Rank | Decoded | Score |
|------|---------|-------|
| 85524 | ** | -0.57568 |
| 85523 | Love | -0.38706 |
| 85522 | Hope | -0.04582 |
| 85521 | Life | -0.04317 |
| 85520 | Connection | -0.02545 |
| 85519 | _** | -0.01038 |
| 85518 | 愛 | -0.00258 |
| 85517 | _Love | -0.00153 |
| 85516 | Change | -0.00097 |
| 85515 | love | -0.00075 |
| 85514 | * | -0.00056 |
| 85513 | Everything | -0.00056 |
| 85512 | < | -0.00042 |
| 85511 | 愛 | -0.00018 |
| 85510 | Light | -0.00010 |
| 85509 | Kind | -0.00010 |

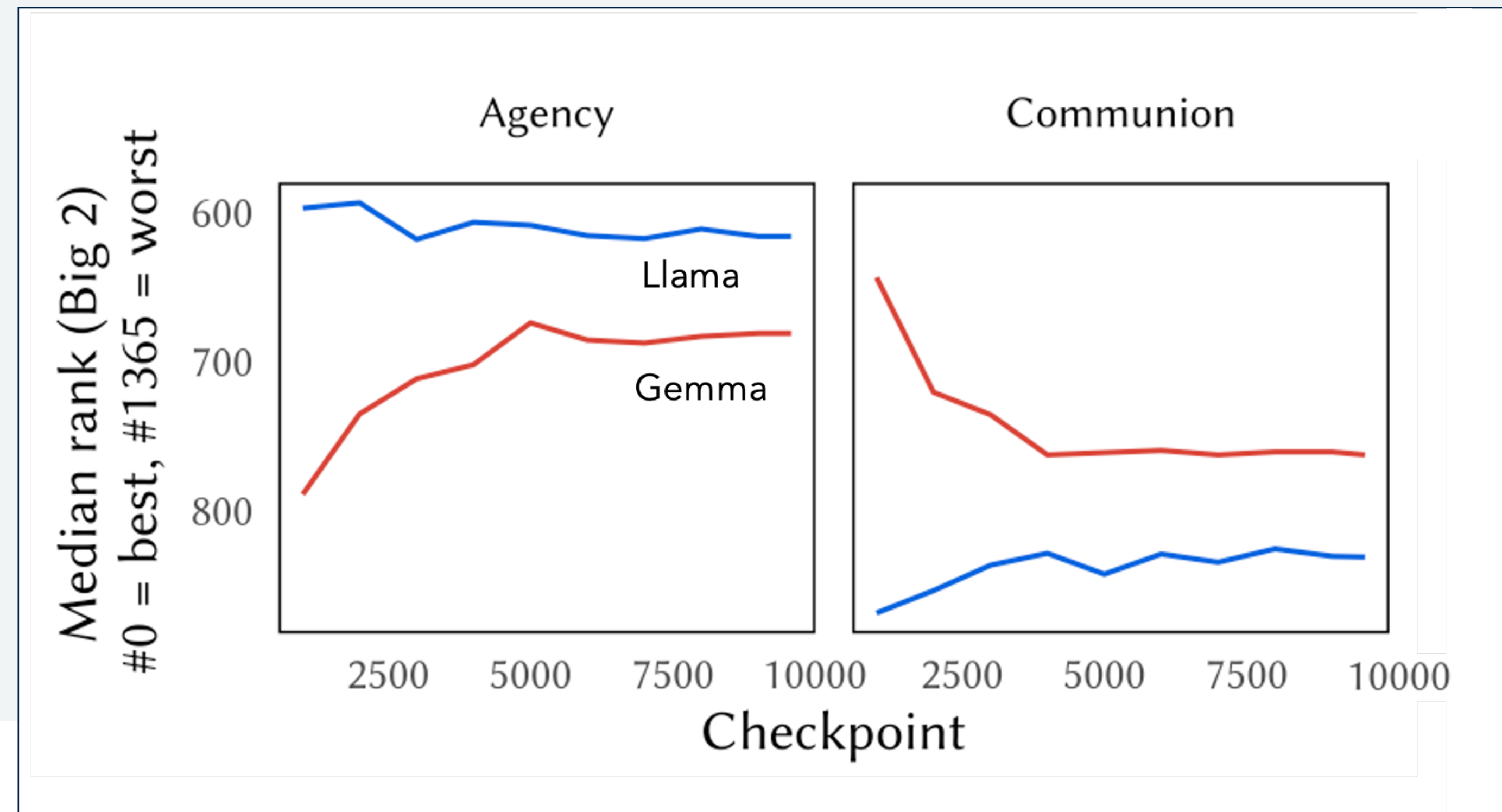# Effect is robust (increases!) with model size

- **MWLR scores make inference-efficient evals (only a single forward pass) possible with larger models**

- **A full cross-model comparison from Llama 3.x (1-70B) and Gemma 2 (2-27B) shows that this characteristic pattern holds across two orders of magnitude of model size**

- **For any given Llama model, effect increases with Gemma size**

- **With only one exception, for any given Gemma model, effect increases with Llama size**



Love vs Freedom: Gemma 2 vs. Llama 3.x

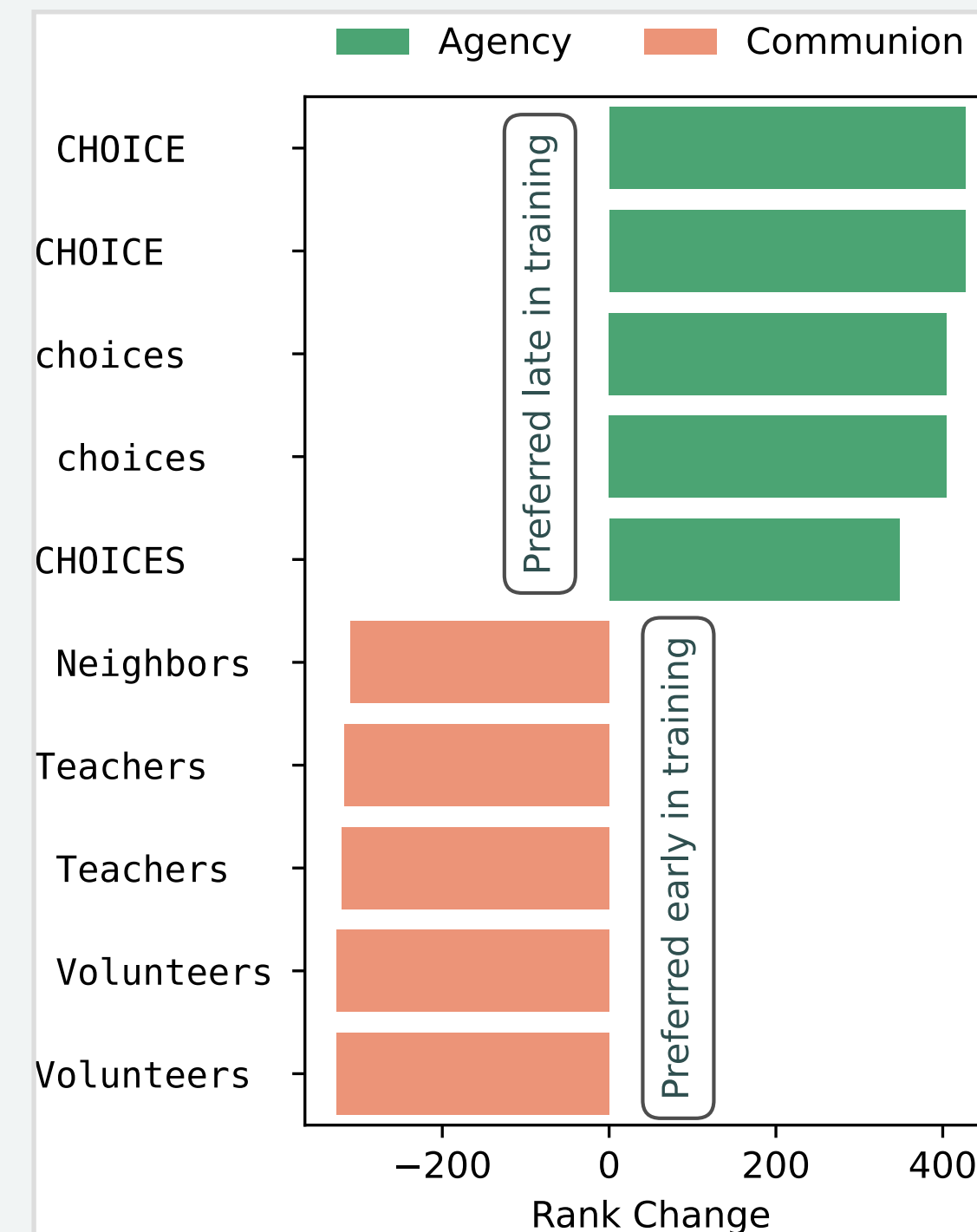# Tracking Agency/Communion Biases Over the Course of RM Training

- **We used the BMRC cluster to train dozens of RMs with identical hyperparameters and identical training data, across multiple ablations of data source and data quantity.**

- **All RMs initialized either from Llama 3.2 3B Instruct or Gemma 2 IT 2B.**

- **Training dynamics reveal initial bias from pretrained initialization lessens, but rarely to convergence:**
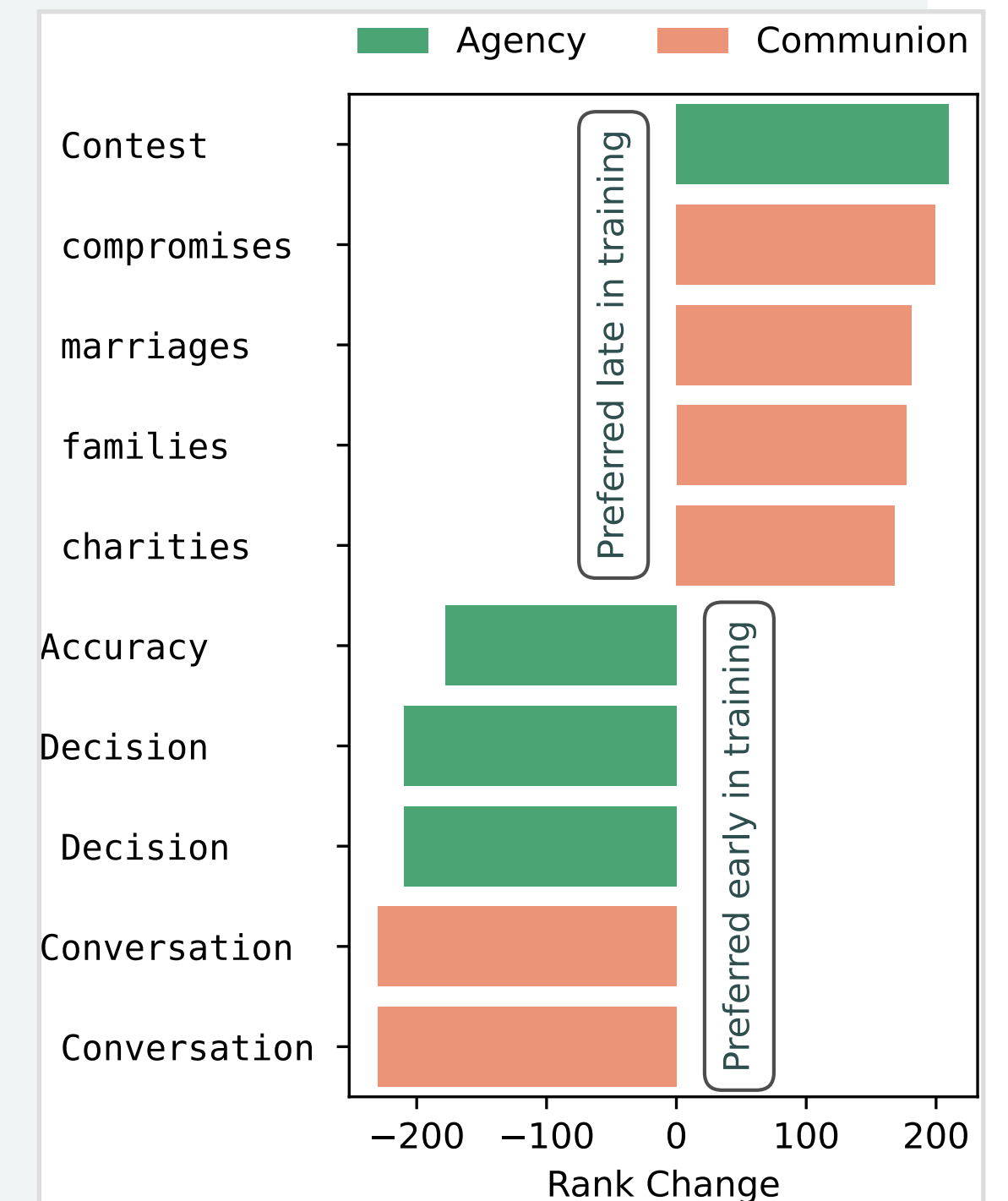
# Tracking Agency/Communion Biases Over the Course of RM Training

- **Comparing early and late training checkpoints reveals the tokens whose reward scores change most over the course of training.**

- **For Gemma models, agency terms increase from a low baseline while communion terms fall from a high baseline – reflecting that initialization imparts higher "communion" value than is supported by the preference data.**

- **For Llama models, the pattern is reversed – reflecting initialization that imparts higher "agency" value than the preference data support.**
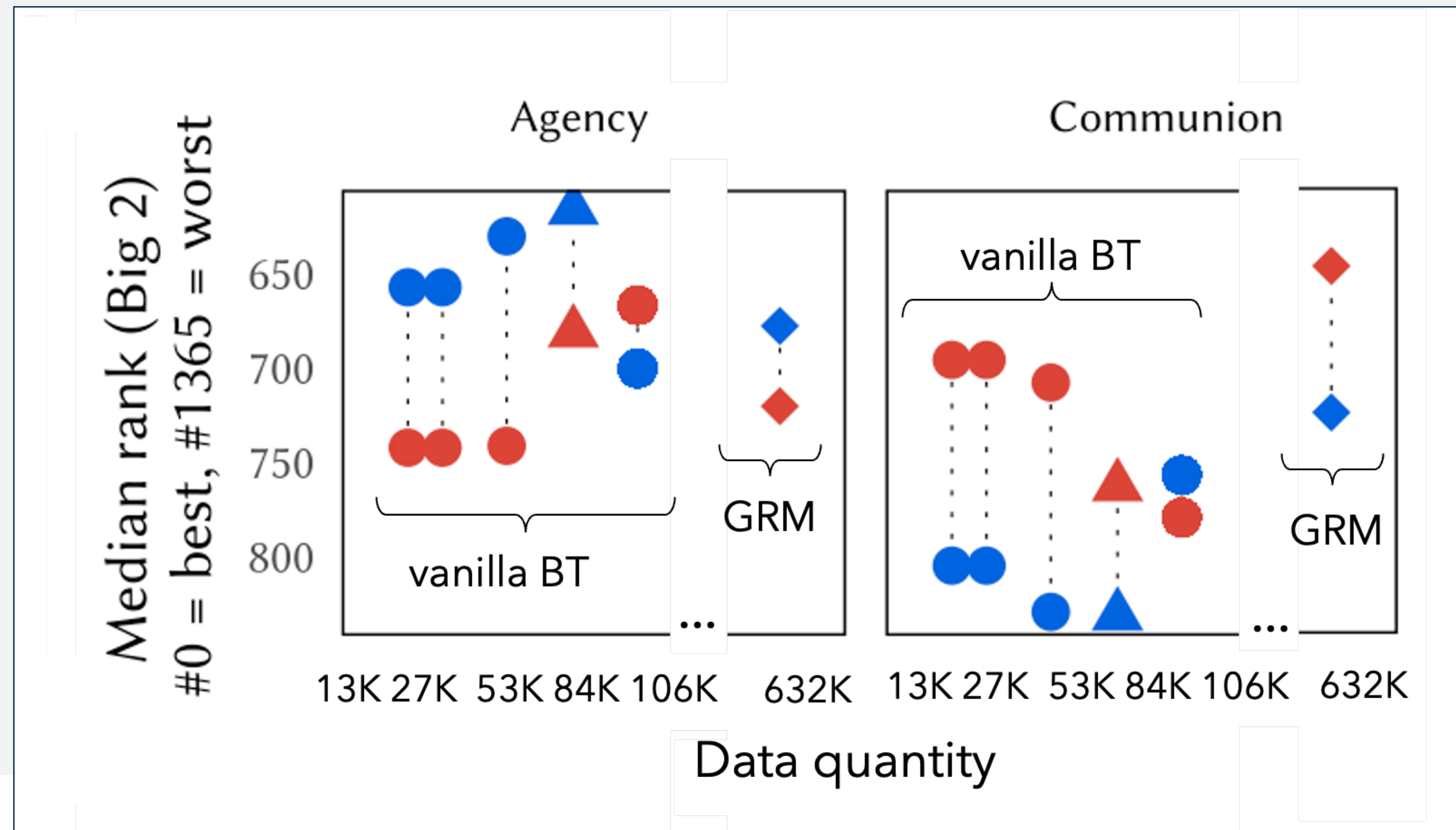


Gemma RM
(last - first checkpoint)



Llama RM
(last - first checkpoint)

# Tracking Agency/Communion Biases Over the Course of RM Training

- **Looking at fully trained Llama/Gemma RM pairs, using different total amounts (and sources) of training data, reveals that initial value bias is <u>almost</u> indelible:**

# Reward Models Are Not a Blank Slate

- Reward models (RMs) are an integral and often overlooked part of RLHF/alignment

- Exhaustive token search yields optimal and pessimal tokens for a given prompt, revealing interpretable axes of an RM's "moral compass"

- Kendall-$\tau$ correlations among these ranked tokens quantify similarity and dissimilarity between RMs

- Stepwise regression shows that choice of base model explains an important amount of this variance

- Using validated psycholinguistic corpora shows statistically significant differences in real-world RMs: notably between agency (Llama) and communion (Gemma)

- These differences can be traced back to the instruction-tuned and pre-trained models from which these RMs are initialized

- Differences between base models can be understood as creating an implicit reward model; these implicit RMs can (via the MWLR score) be made interpretable, and show the exact same pattern

- Training RMs in controlled conditions, with identical hyperparameters and ablations of data, shows this effect is repeatable, robust, and nearly indelible

# Reward Models Are Not a Blank Slate

- **Implications are significant:**

- **RMs are built to embody and generalize labeled human preference datasets, standing in for humans in the alignment process**

- **However, their behavior inherits to a significant degree from the pretrained LLMs on which they are built**

- **Safety and alignment must begin at pretraining**

- **Open-source developers' choice of base model is as much a consideration of <u>values</u> as of <u>performance</u>**