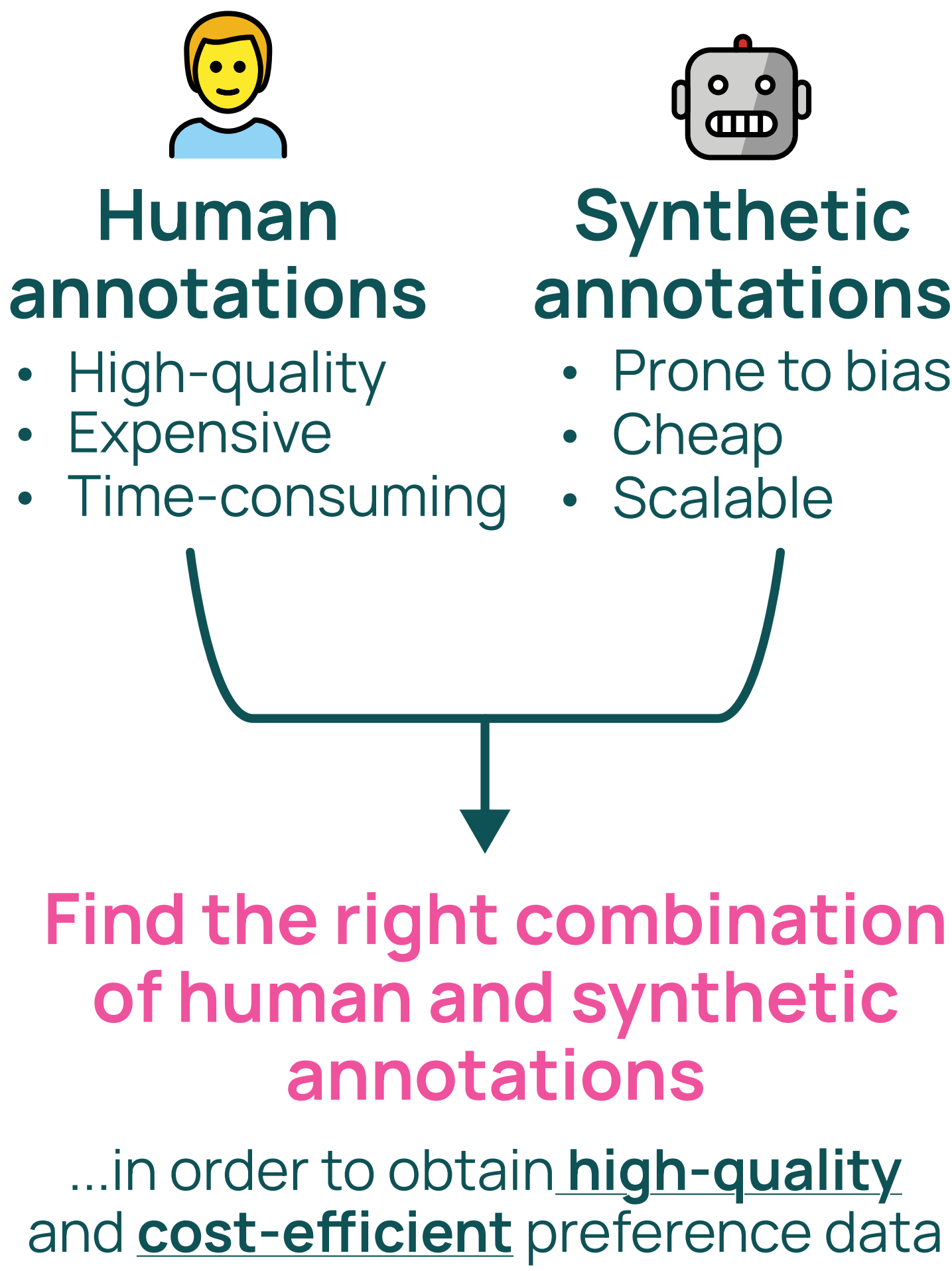


Hybrid Preferences: Learning to Route Instances for Human vs. AI Feedback

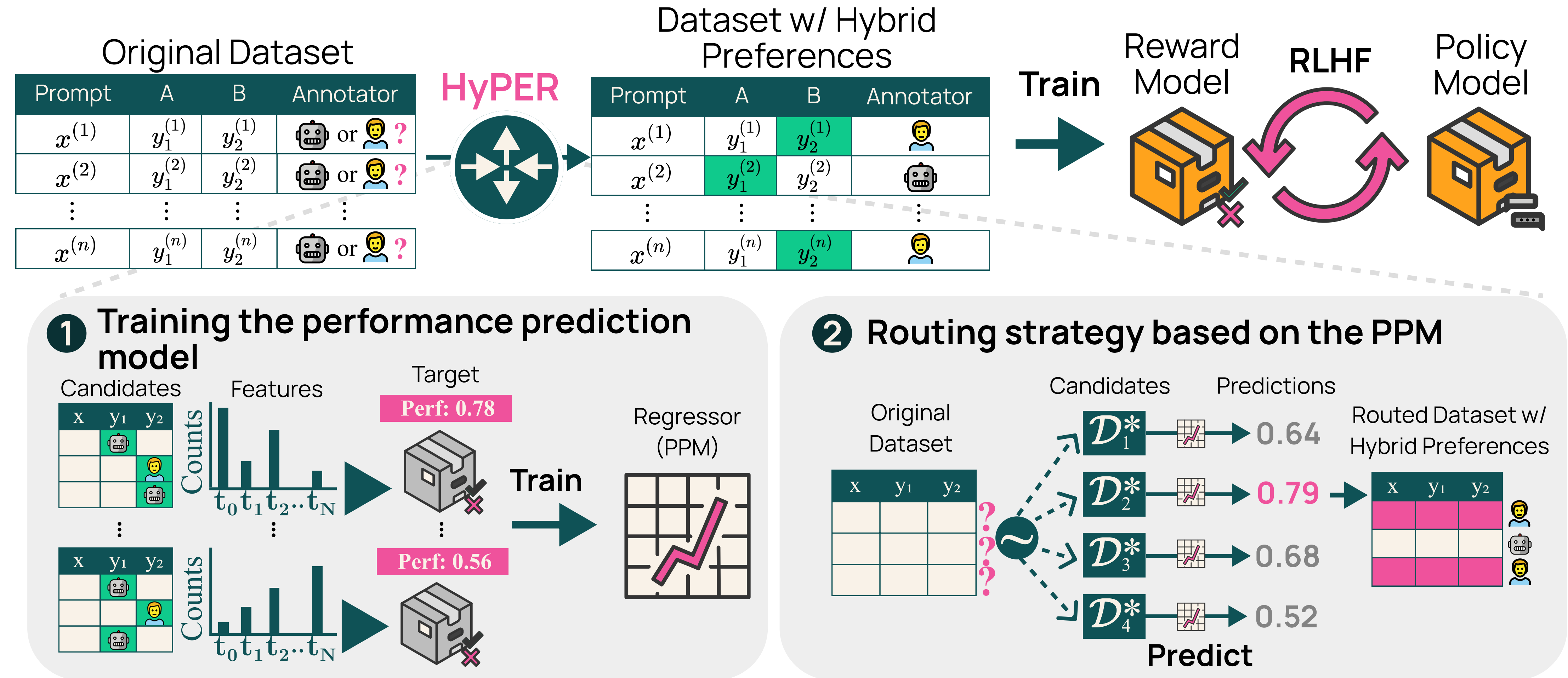


L.J.V. Miranda, Y. Wang, Y. Elazar, S. Kumar, V. Pyatkin, F. Brahman, N. A. Smith, H. Hajishirzi, P. Dasigi

Motivation



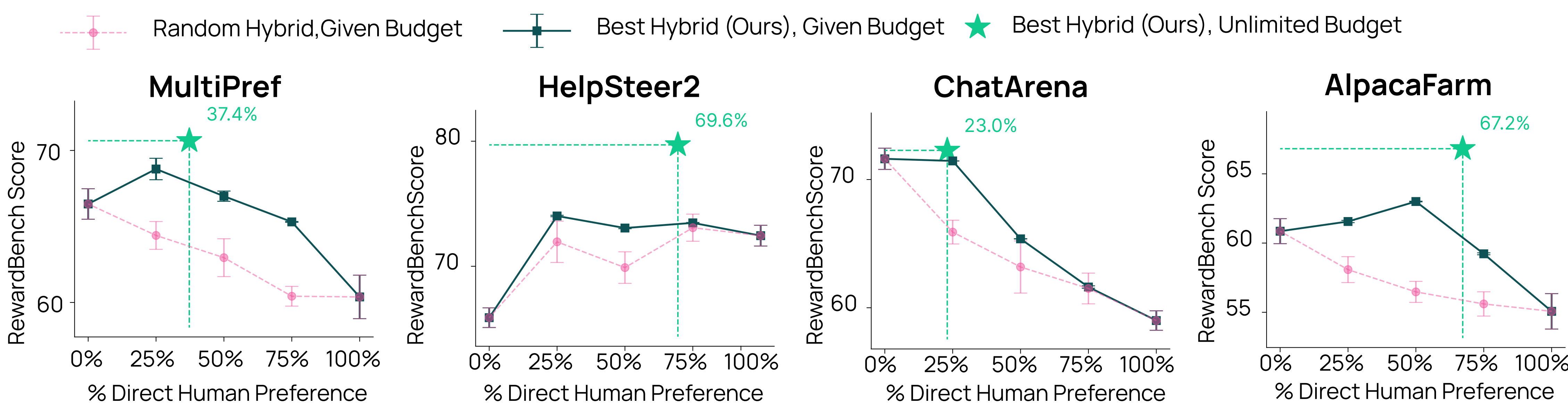
Methodology



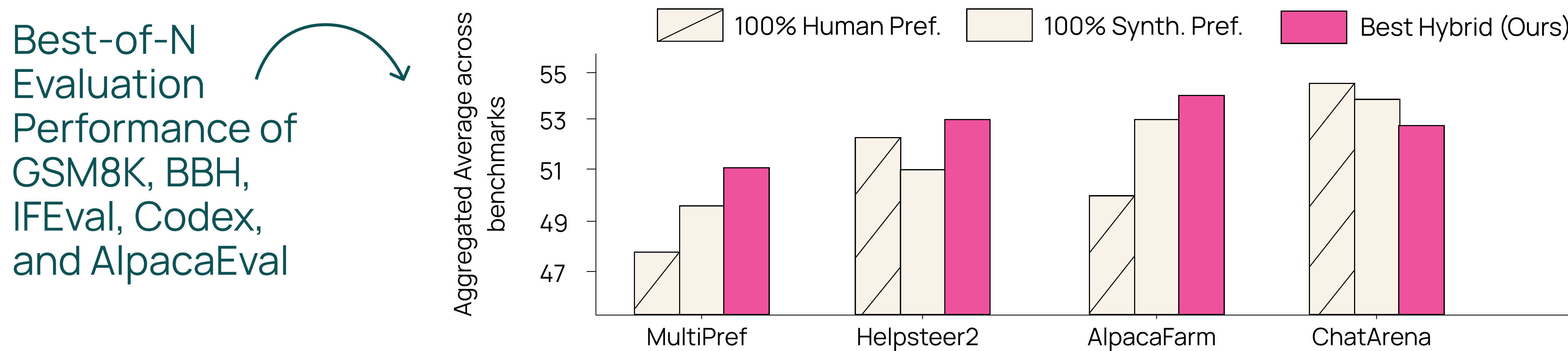
Results

Reward models trained on our routed datasets perform better vs. random / 100% human / 100% LLM on **unseen datasets, benchmarks, and base models.**

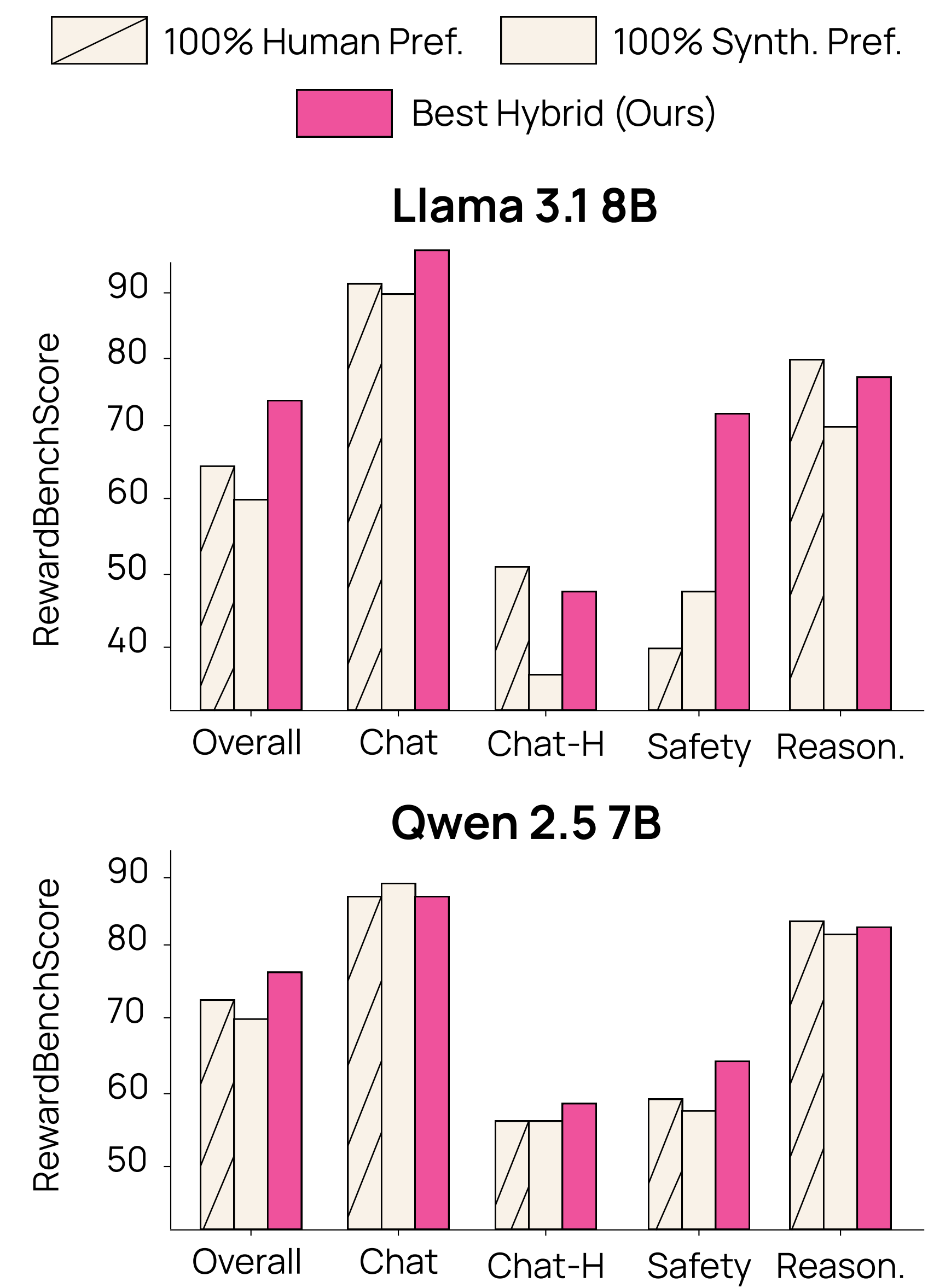
Unseen Preference Datasets



Other Benchmarks

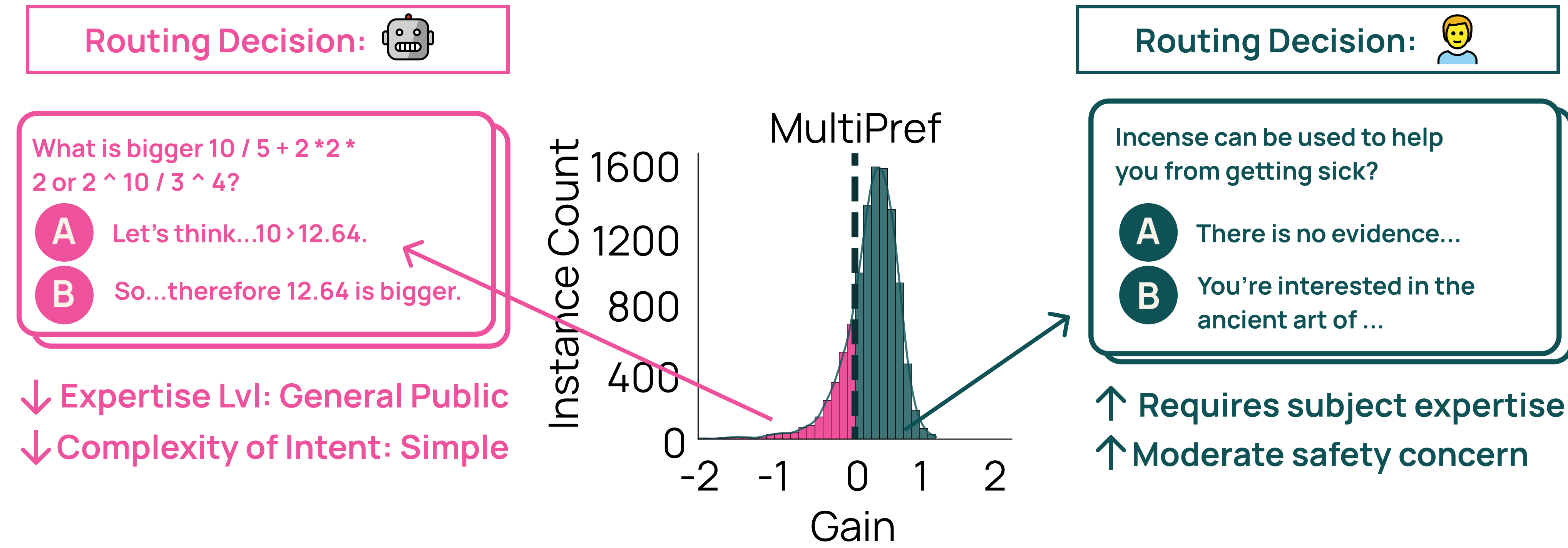


Other Base Models



Analysis

By looking at the features learned by our PPM, we can understand the characteristics of instances suited for human annotation



Contributions

- We find that some preference instances are **better suited to be annotated by humans** than language models.
- We used this to build a **routing framework** for preference data.
- We obtain **fine-grained understanding** of what type of instances benefit human annotations.