# Supplementary Material

## Evaluating the Effectiveness of Error Bars and Multilevel Confidence Bands in Conformalized Regression

M. Dossi[1], J. T. Lin[1], A. M. Smits[2], A. Chatzimparmpas[1]

[1] Department of Information and Computing Sciences, Utrecht University, The Netherlands
[2] Department of Human Centered Technologies, Hanze University of Applied Sciences, The Netherlands

# 1. House Visualization
## Survey Interface and Results from Participant Estimates

House visualizations are presented by complexity level.

For each house, the first slide shows the uncertainty visualizations used in the survey interface, and the second slide presents the corresponding analysis results.
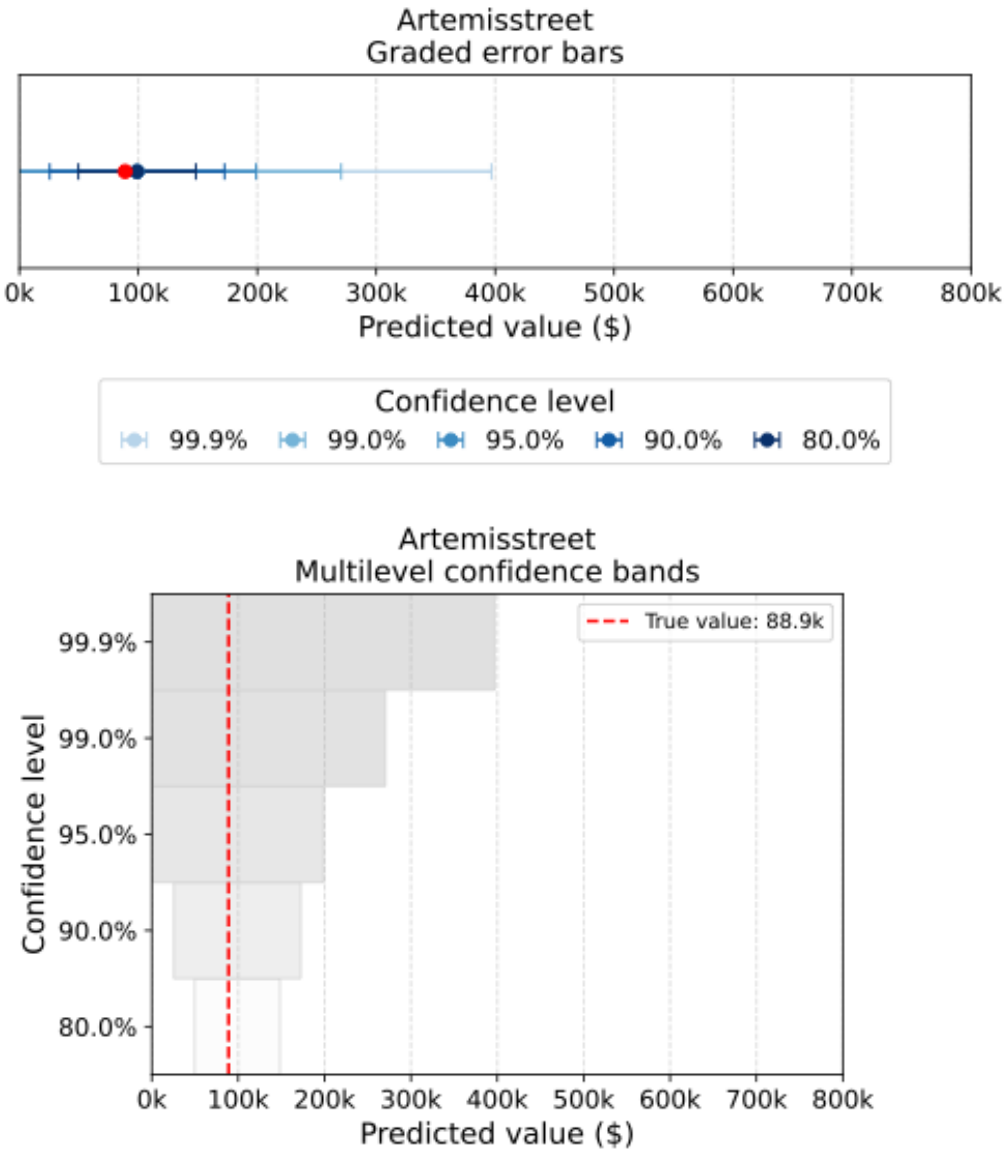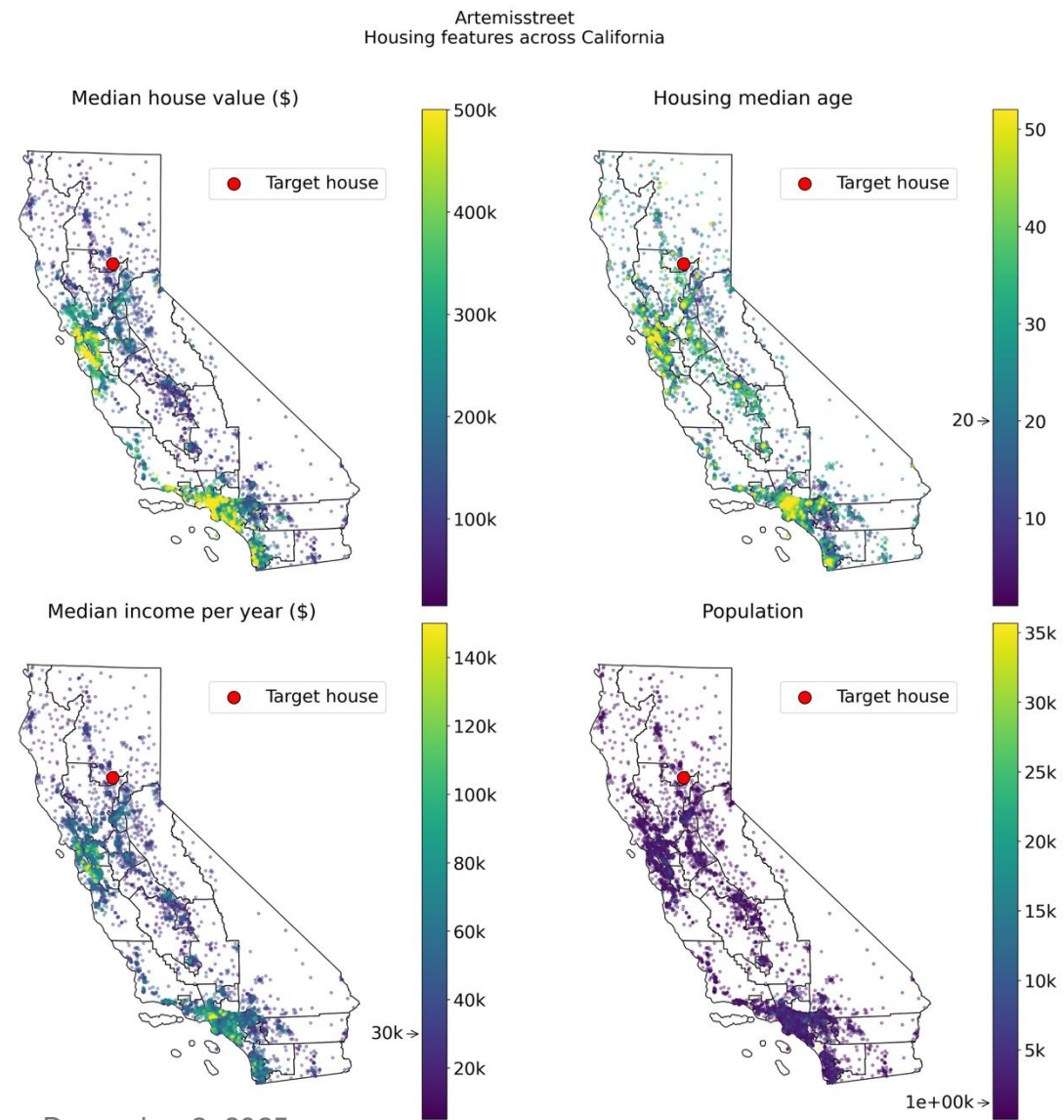
Note that, unlike in the survey, these slides include a red dot (for graded error bars) or a red dashed line (for multilevel confidence bands) to indicate the model's prediction for the target house. This cue was intentionally omitted in the survey.
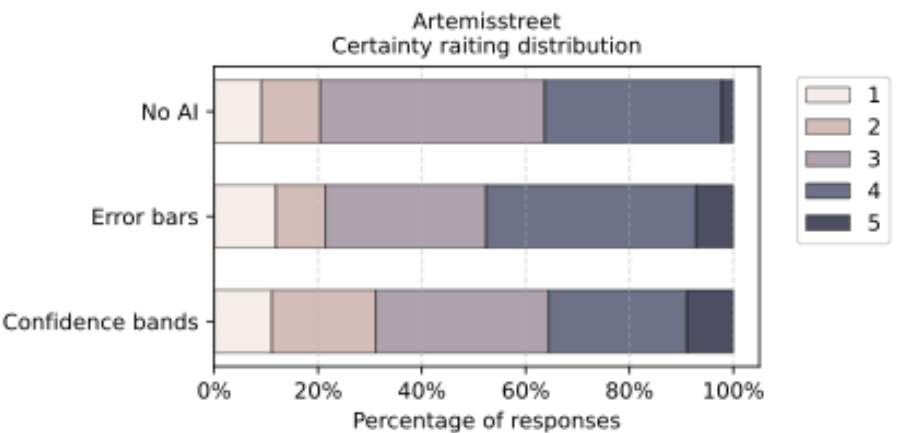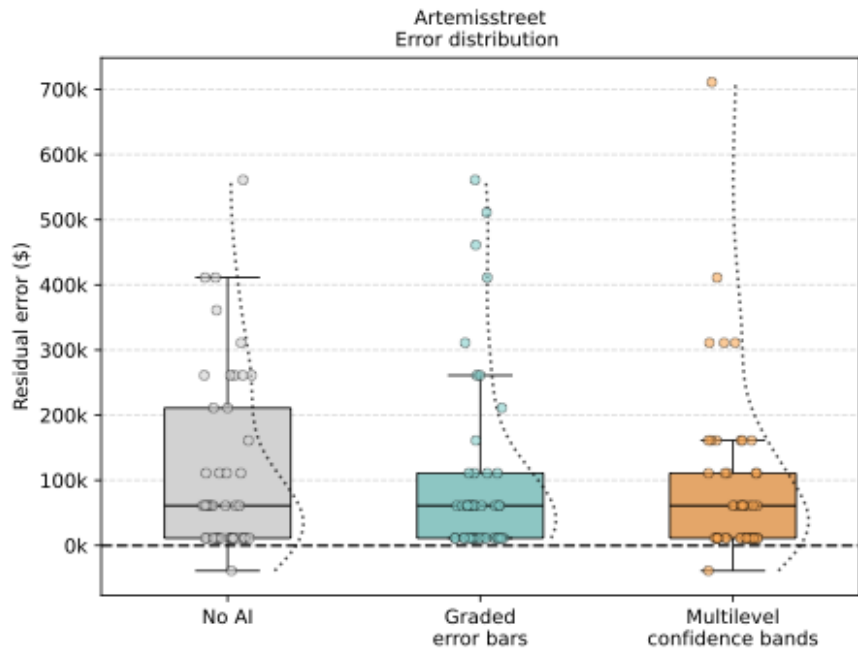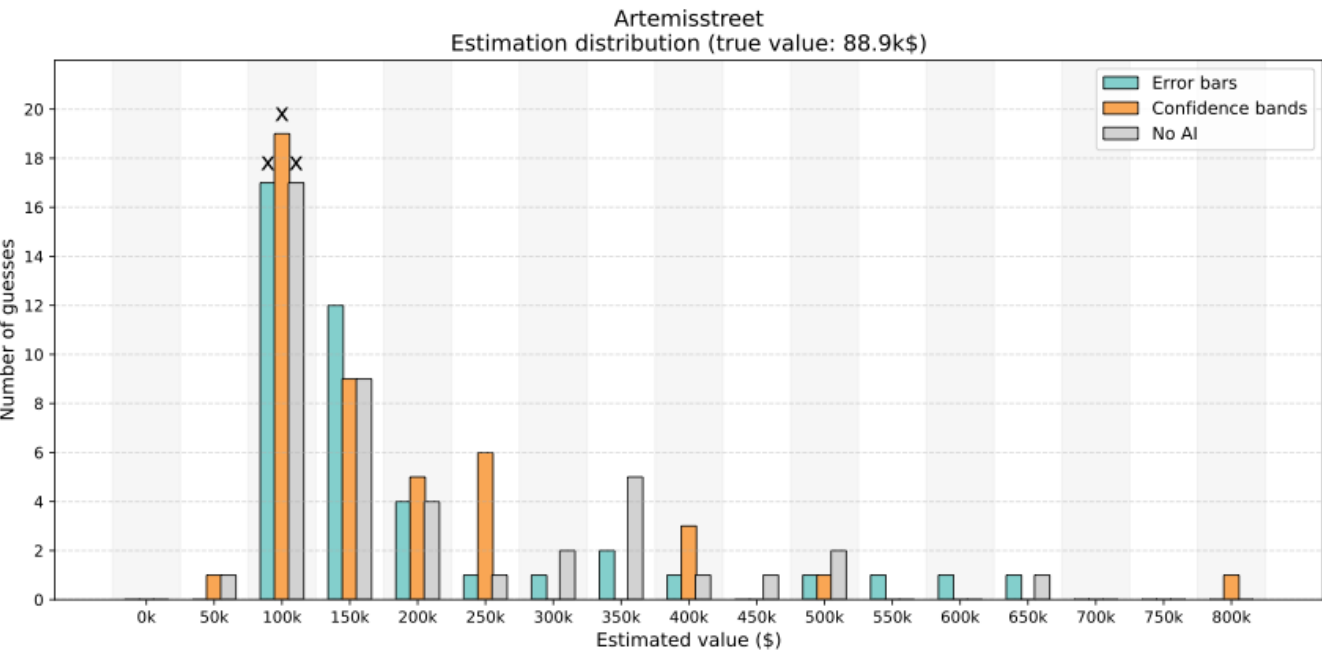
# Task complexity: easy

- 1) Artemisstreet
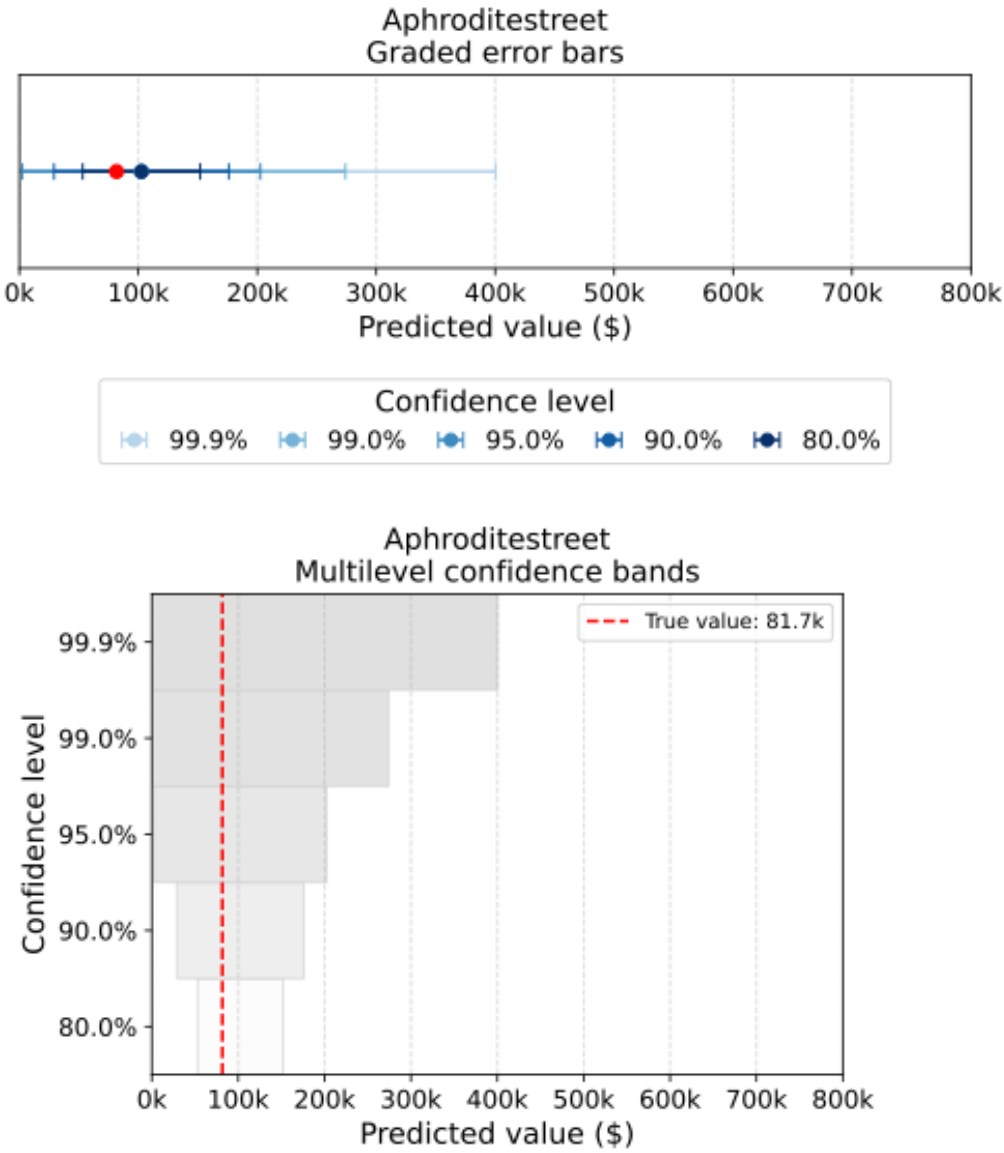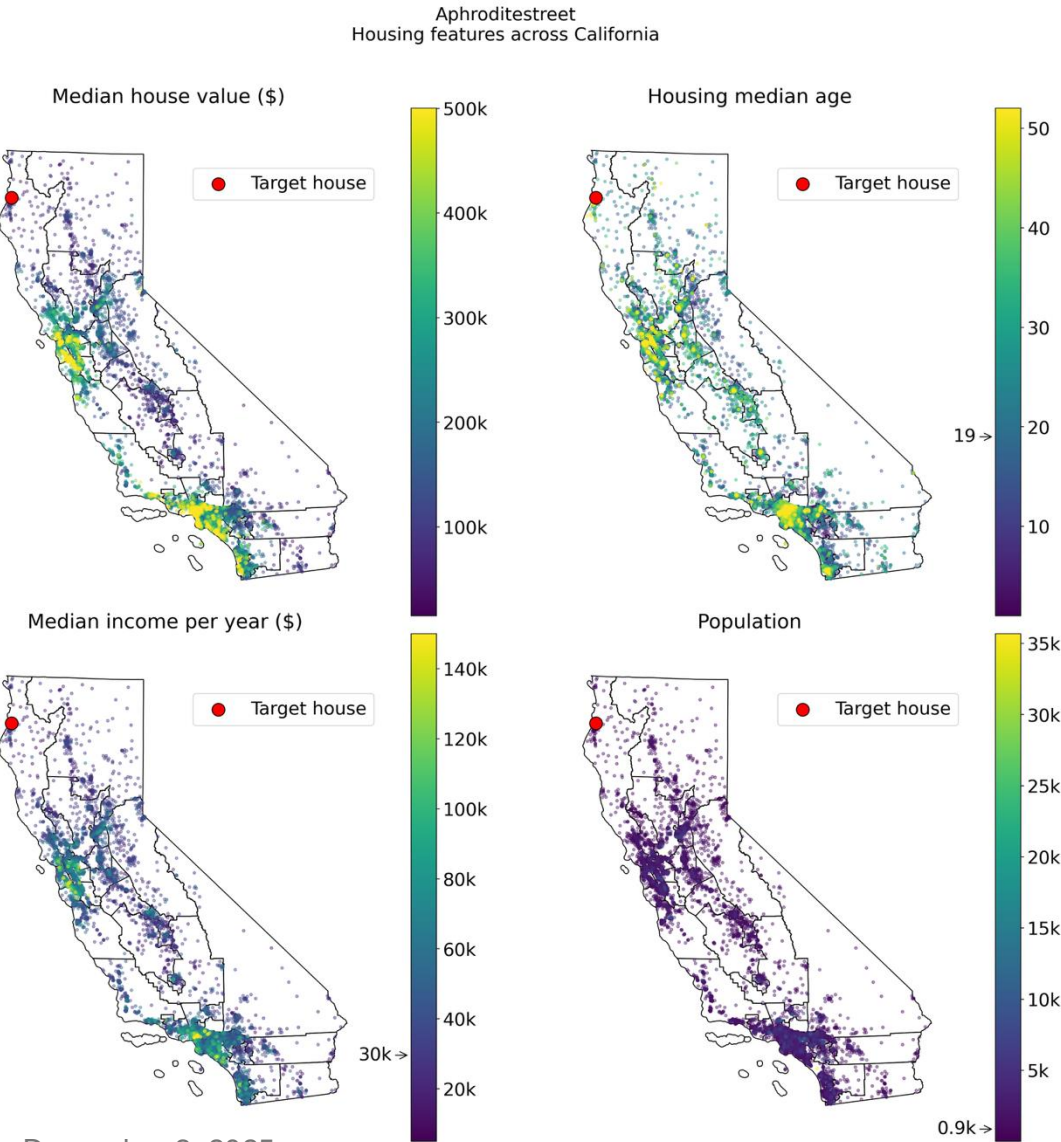- 2) Aphroditestreet
- 3) Apollostreet

# 1) Low error – Artemisstreet

# 1) Low error – Artemisstreet

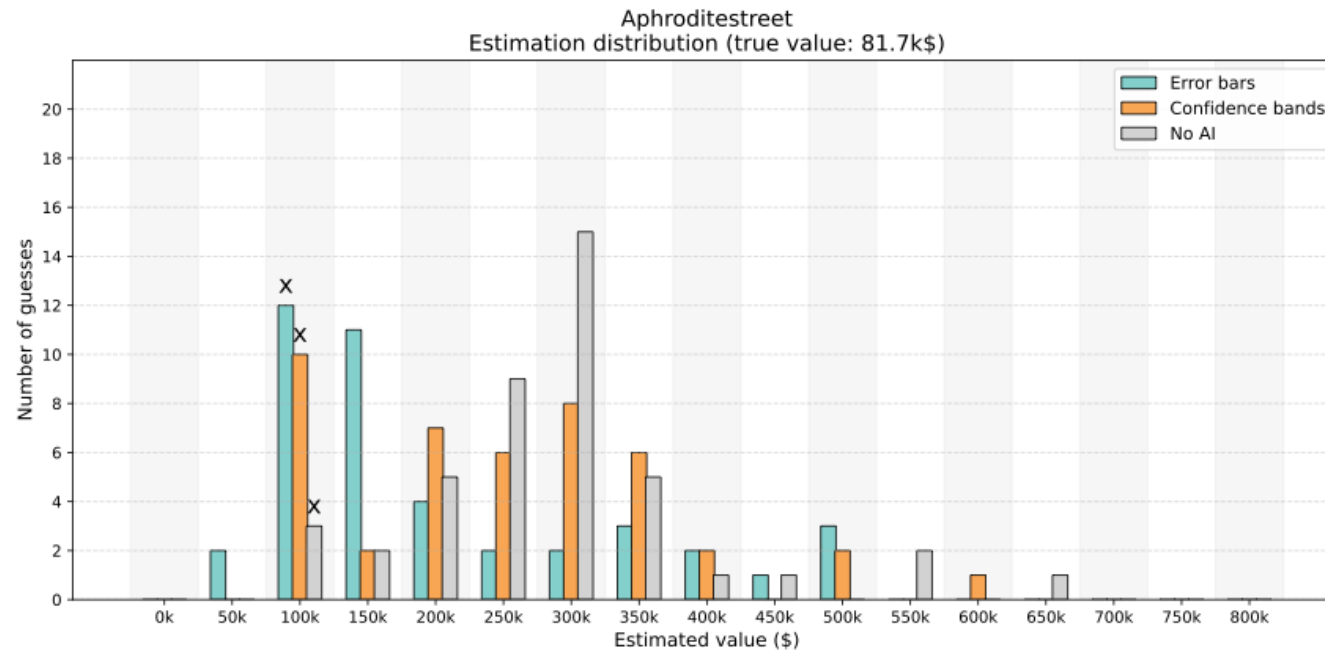Estimates closest to the true value are
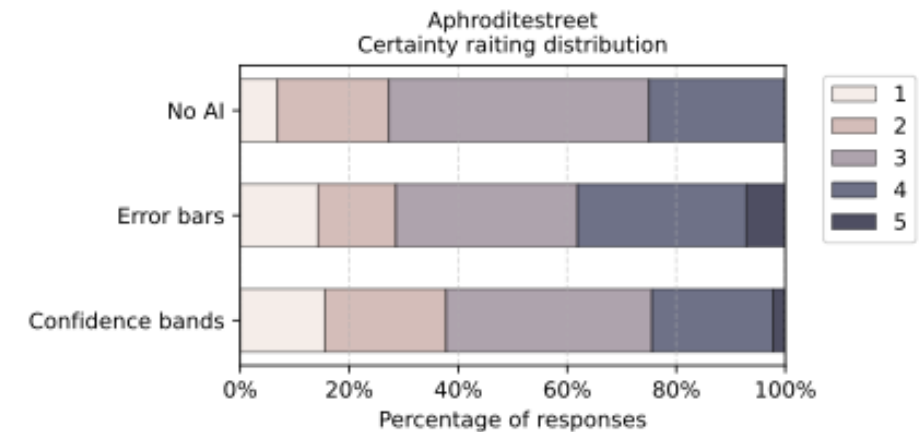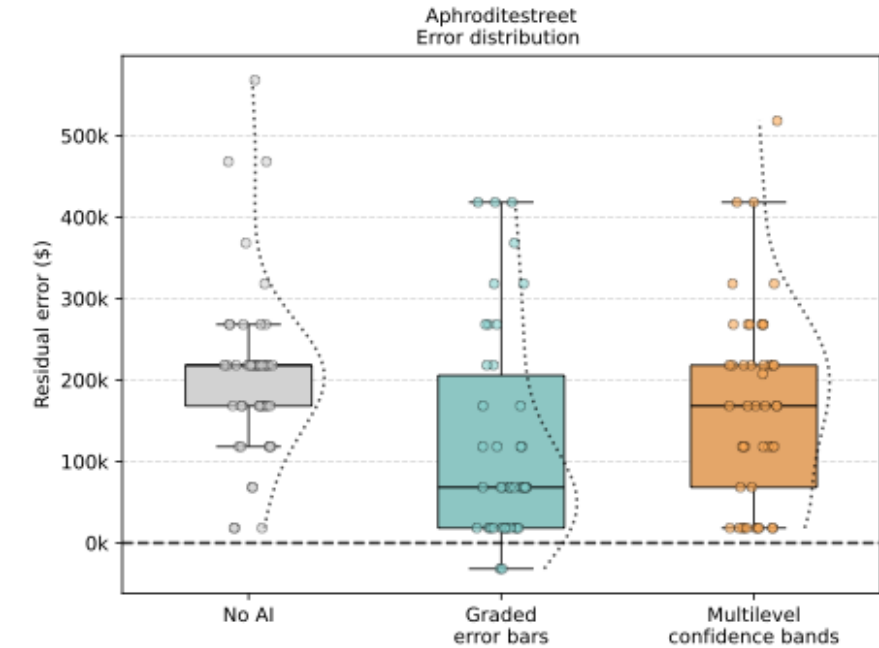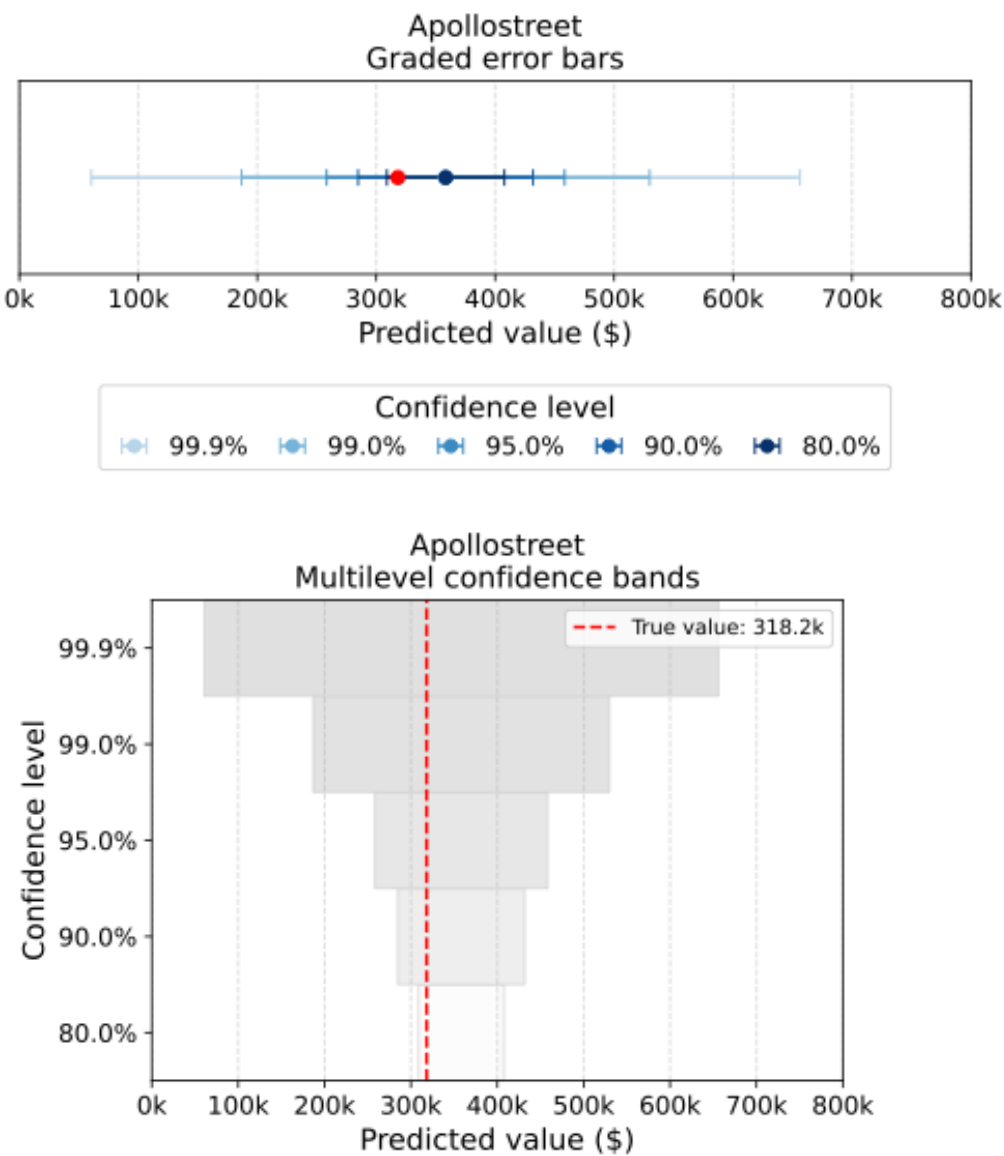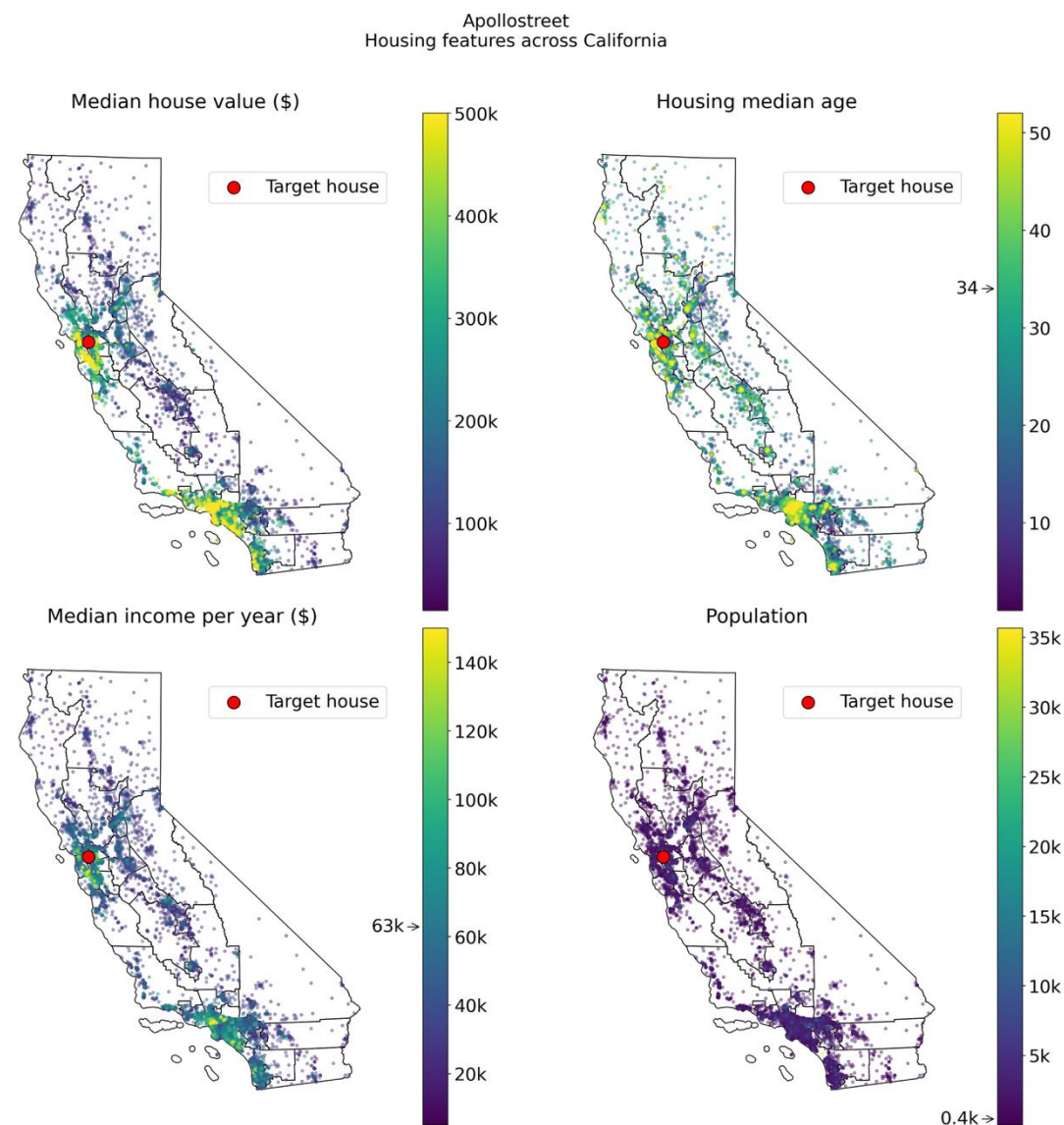marked with an "x" above the bars.

# 2) Low error – Aphroditestreet



Aphroditestreet
Estimation distribution (true value: 81.7k$)



Aphroditestreet
Error distribution



Aphroditestreet
Certainty raiting distribution

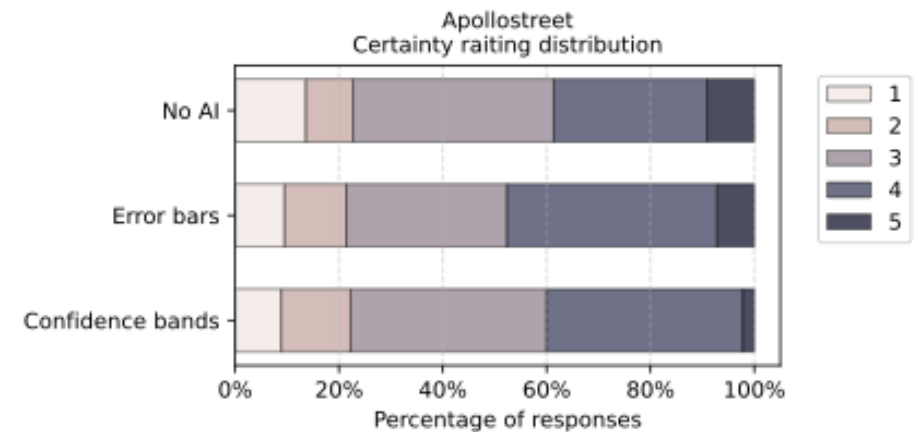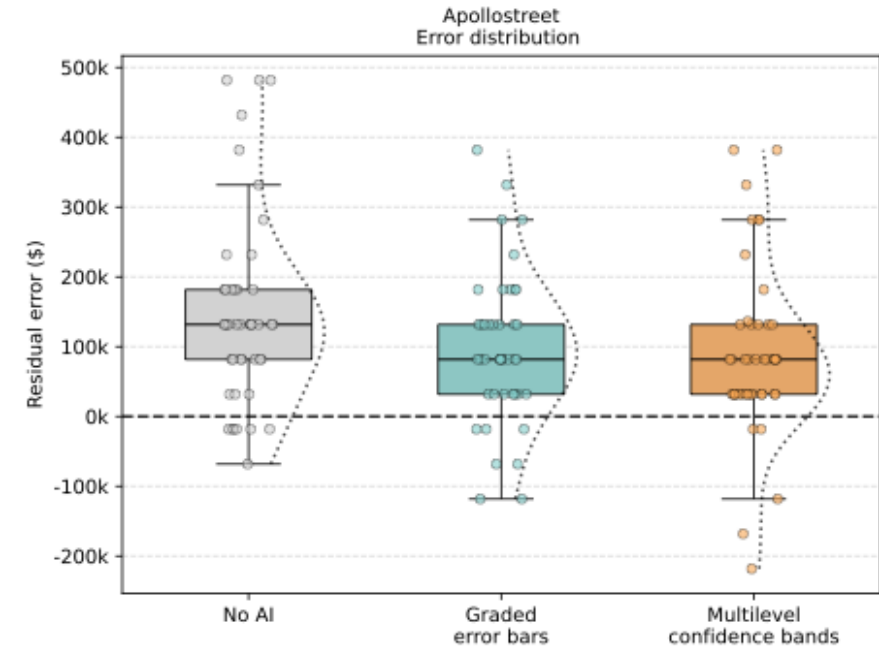Estimates closest to the true value are marked with an "x" above the bars.

Apollostreet
Housing features across California

# 3) Low error – Apollostreet

Results



Apollostreet
Estimation distribution (true value: 318.2k$)

Estimates closest to the true value are marked with an "x" above the bars.



Apollostreet
Error distribution



Apollostreet
Certainty raiting distribution

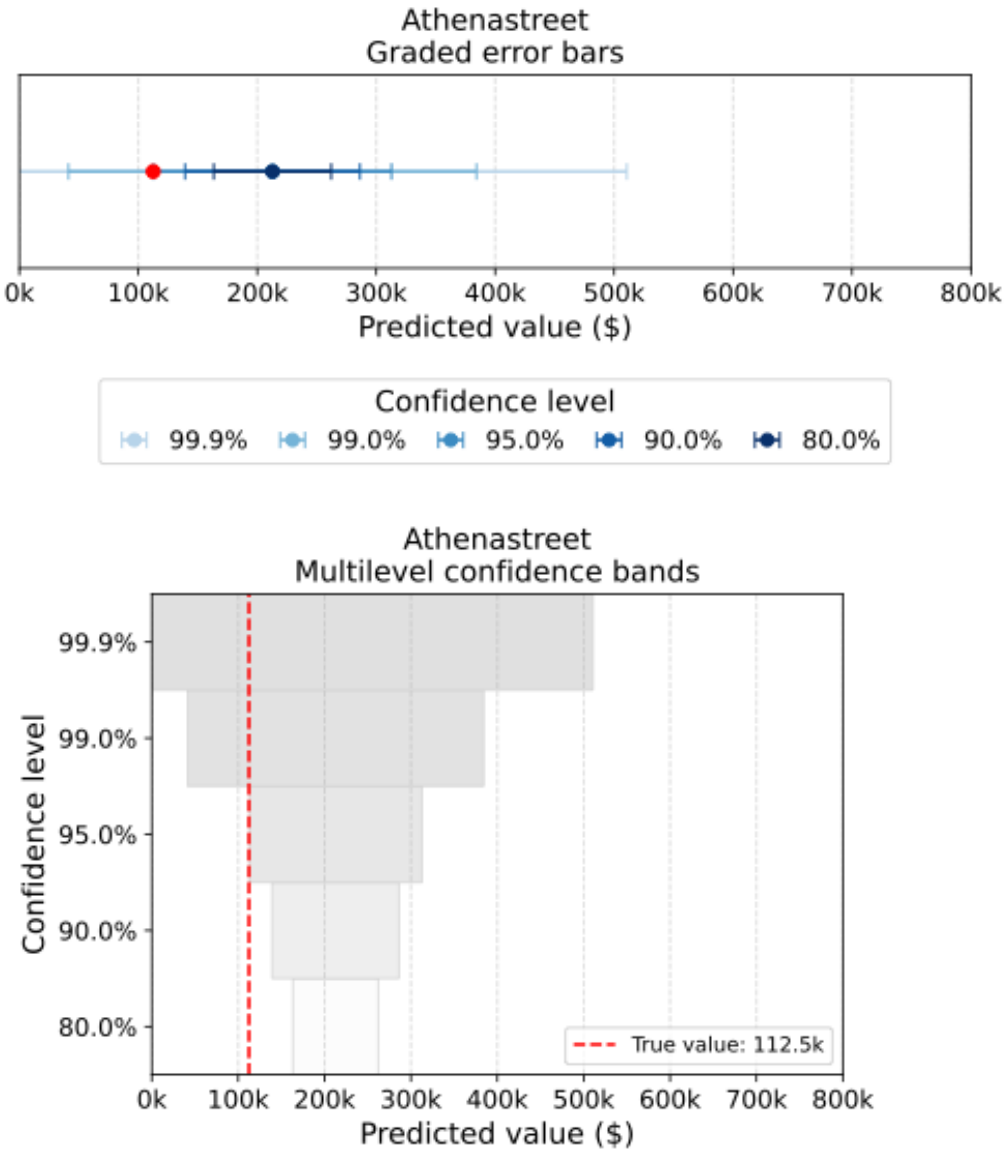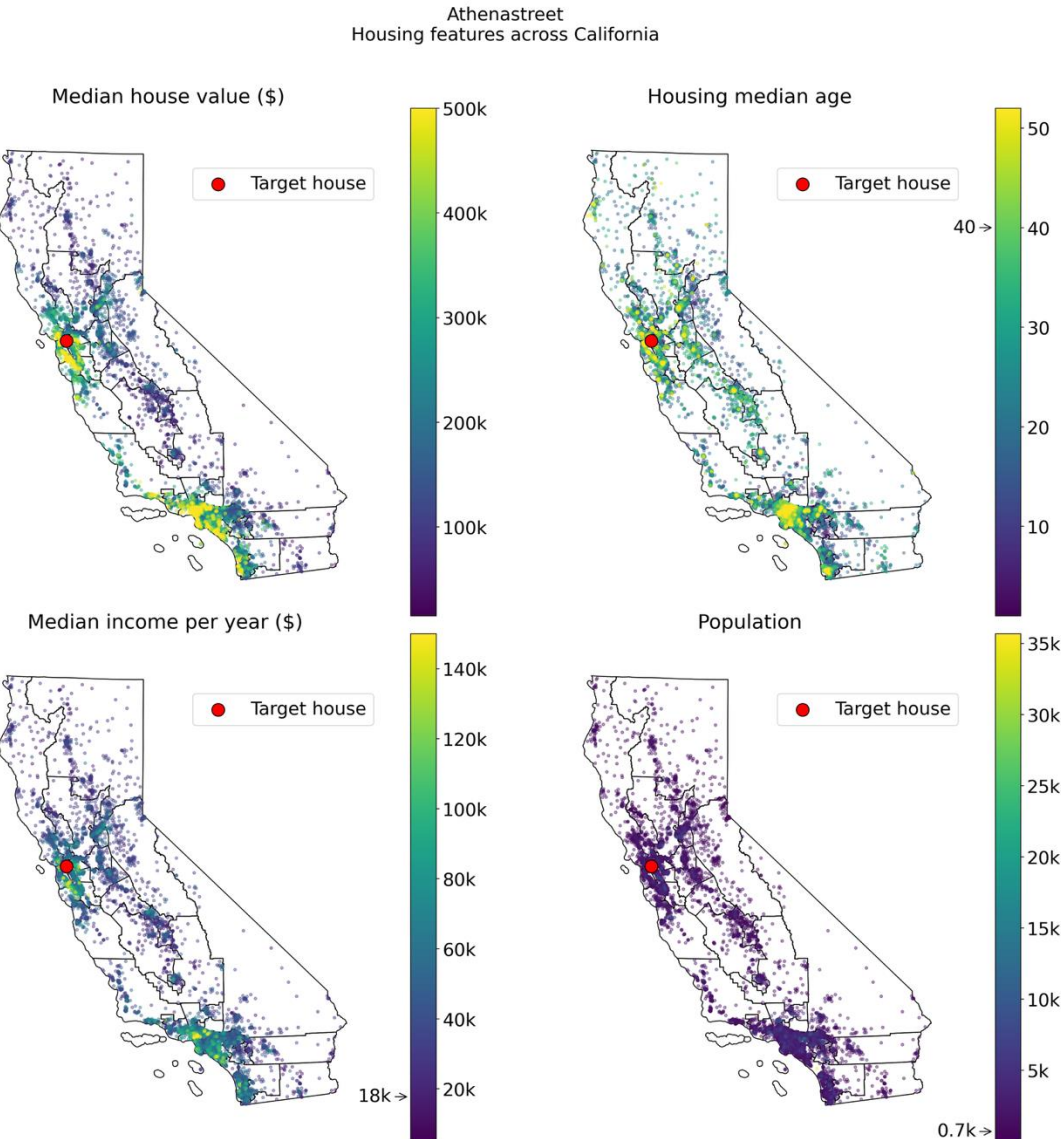December 3, 2025                                                                                    9

# Task complexity: medium

🏠 4) Athenastreet

🏠 5) Poseidonstreet
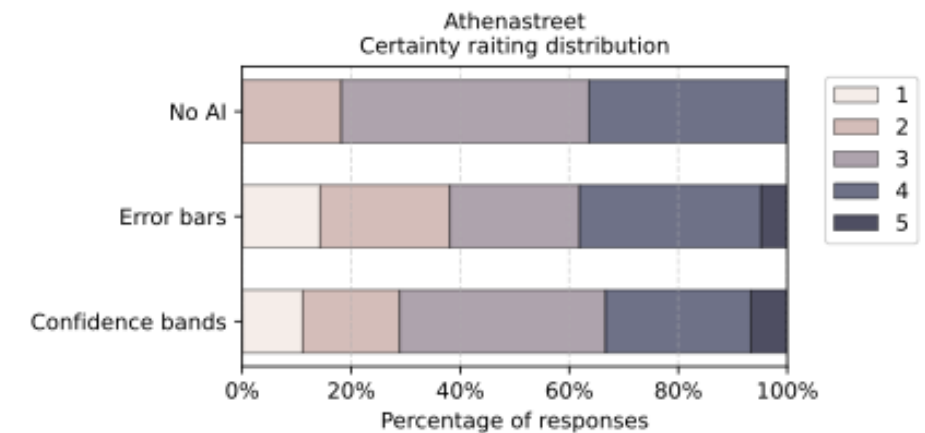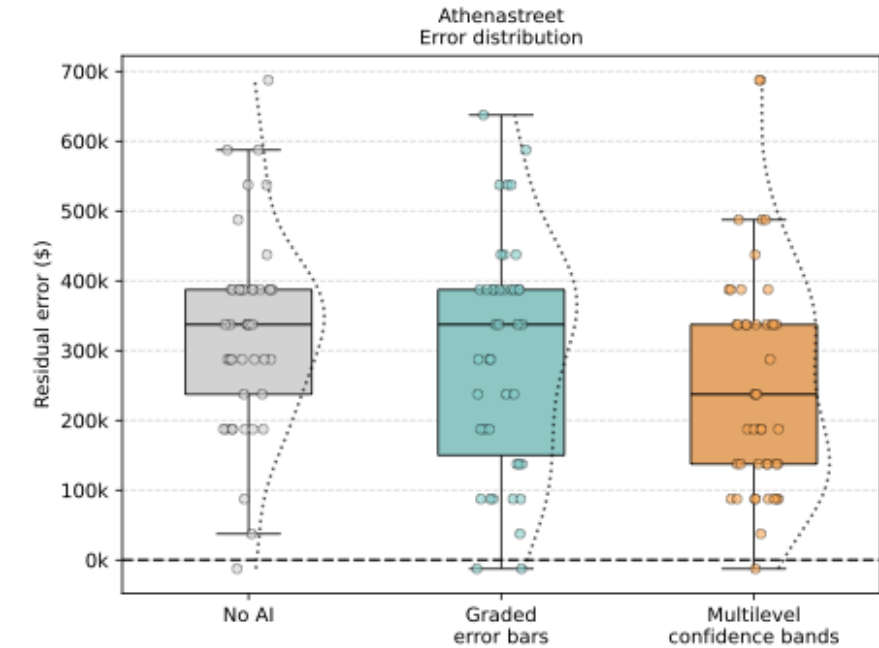
🏠 6) Aresstreet

# 4) Medium error – Athenastreet



Athenastreet
Housing features across California

Median house value ($)

Housing median age

Median income per year ($)

Population

Athenastreet
Graded error bars

Confidence level
99.9%   99.0%   95.0%   90.0%   80.0%

Athenastreet
Multilevel confidence bands
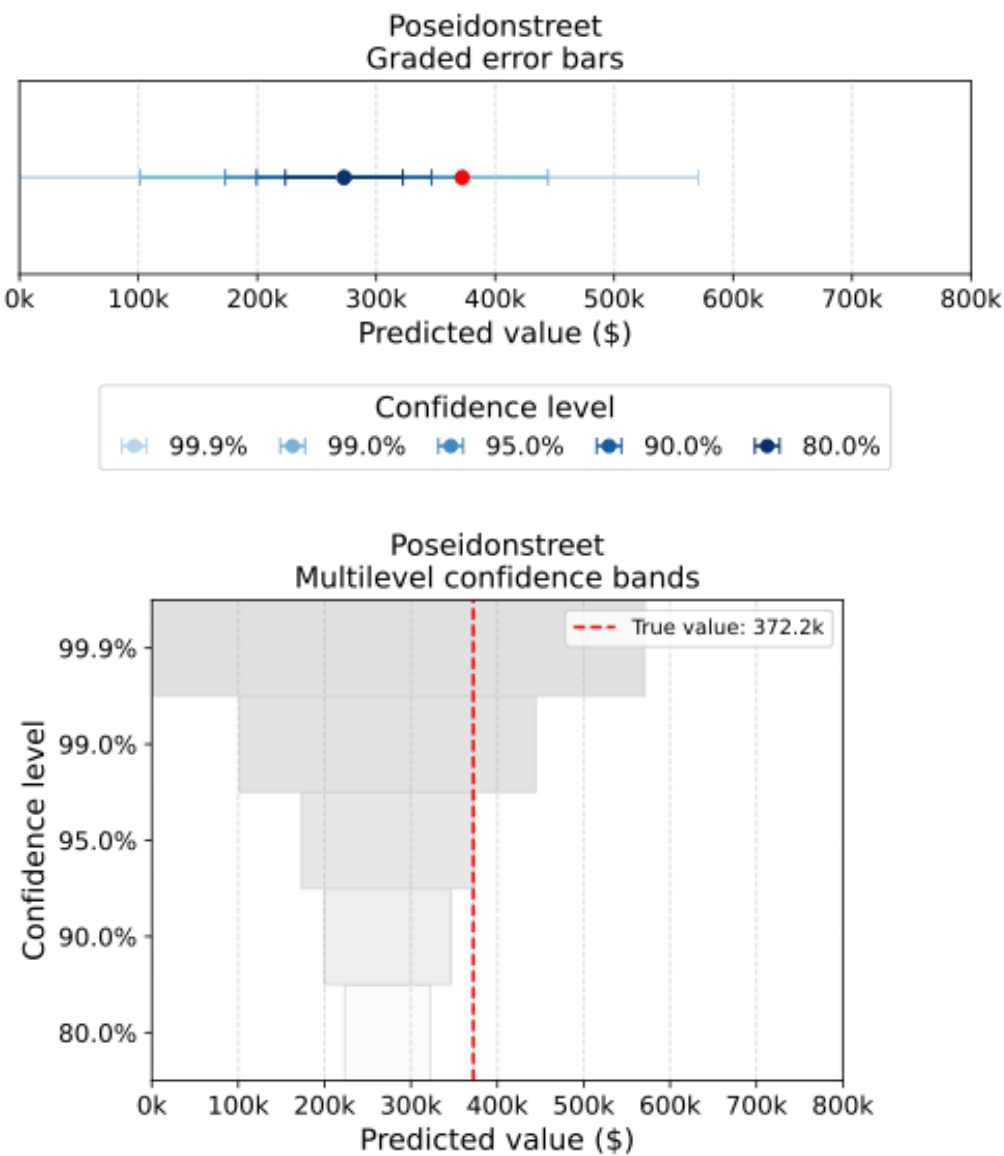
True value: 112.5k

# 4) Medium error – Athenastreet



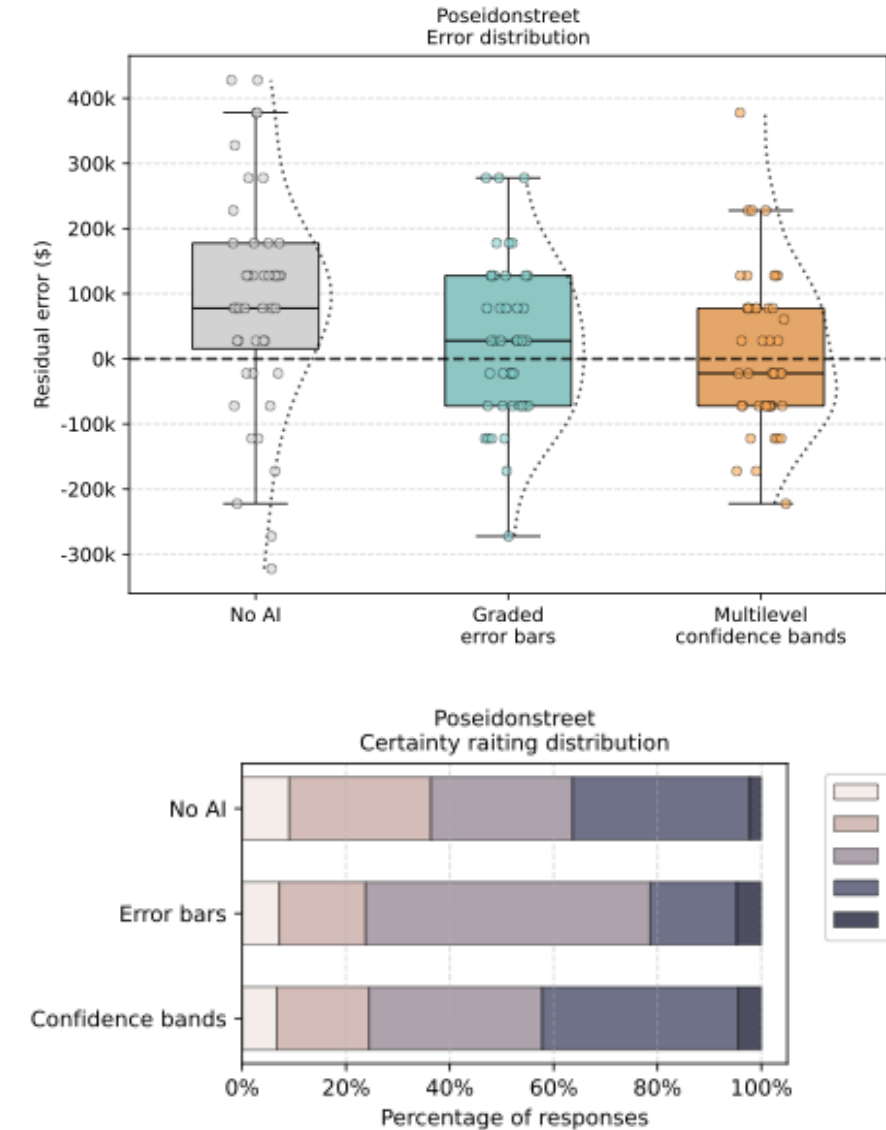Estimates closest to the true value are
marked with an "x" above the bars.

# 5) Medium error – Poseidonstreet
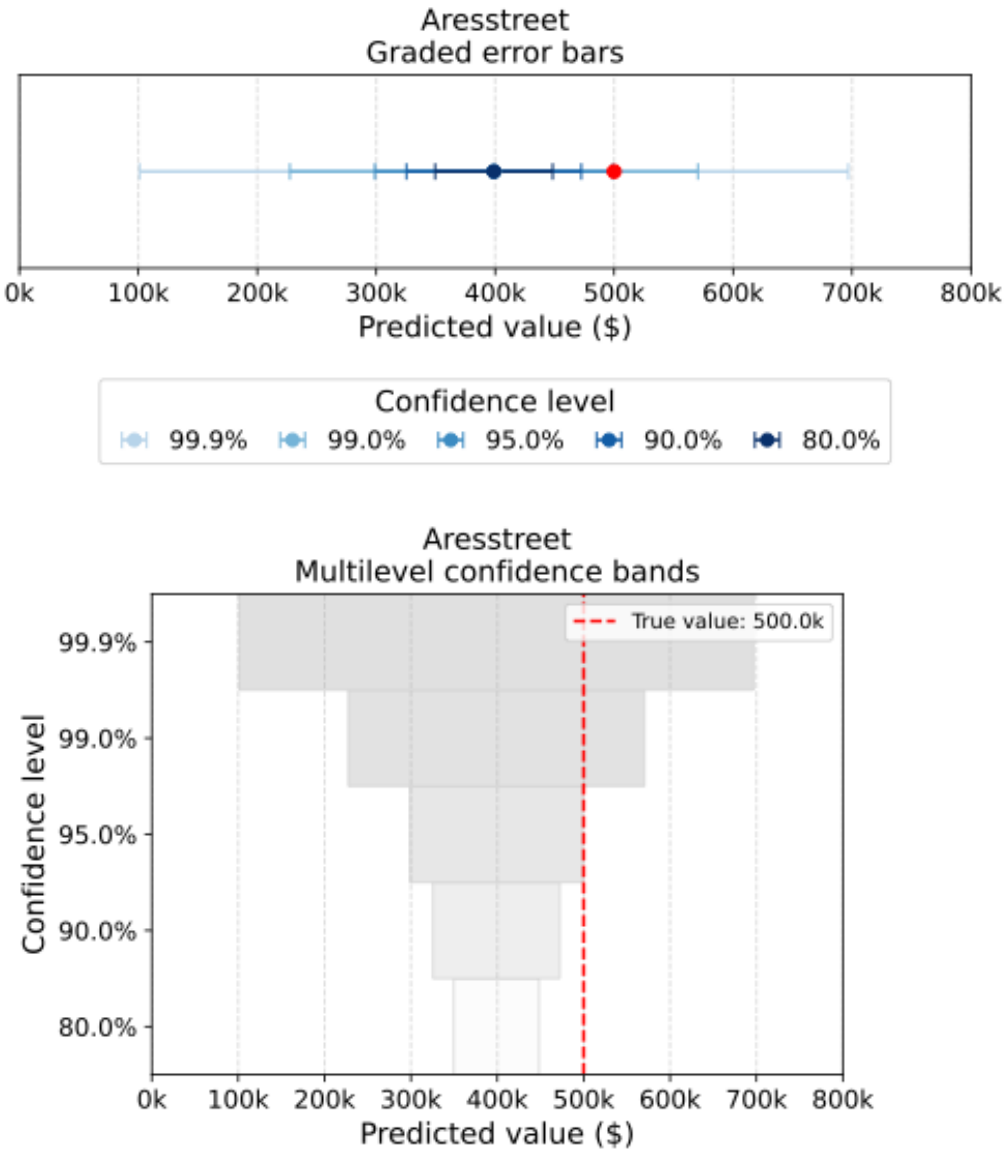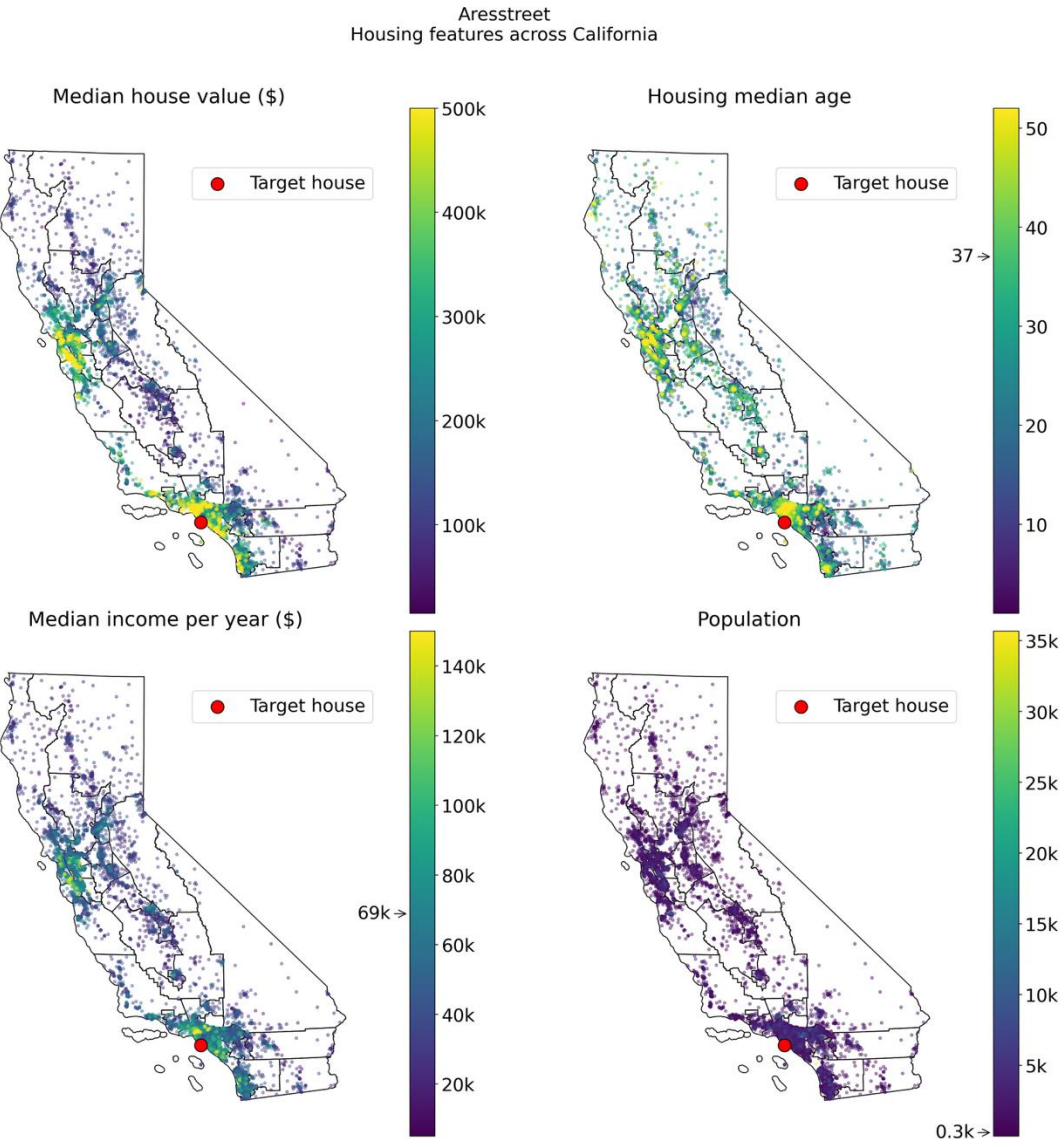
# 5) Medium error – Poseidonstreet



Poseidonstreet
Estimation distribution (true value: 372.2k$)

Estimates closest to the true value are
marked with an "x" above the bars.



Poseidonstreet
Error distribution



Poseidonstreet
Certainty raiting distribution

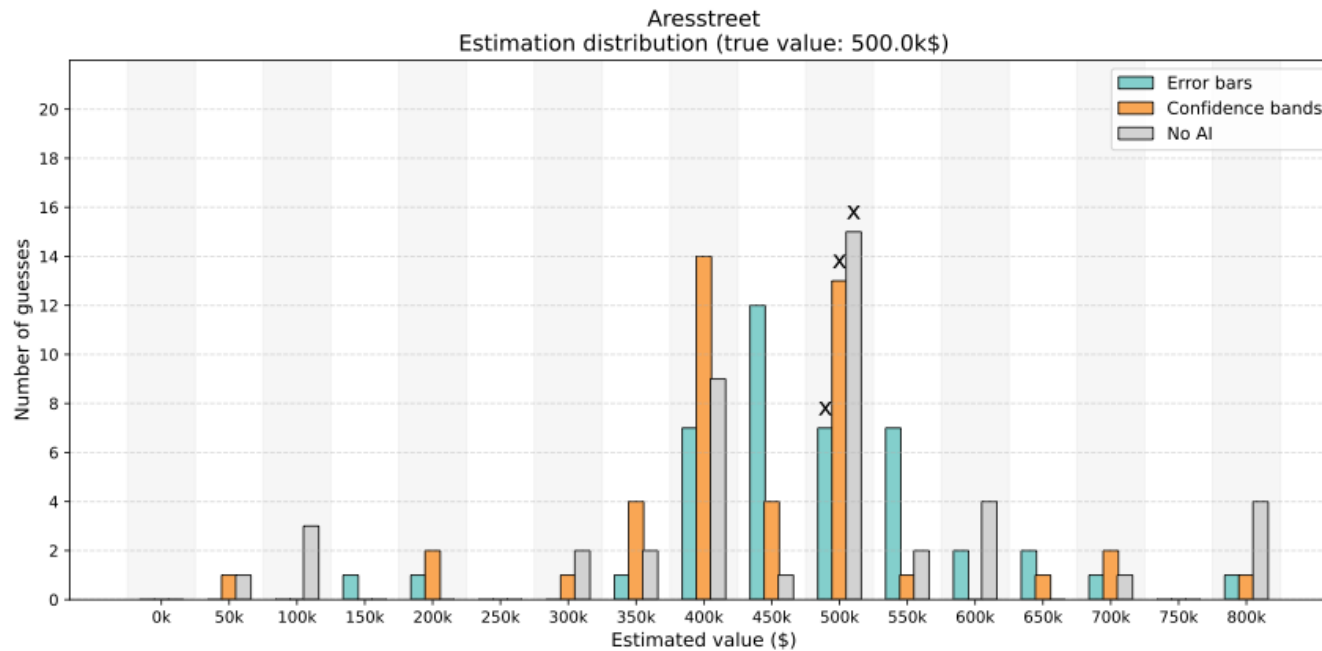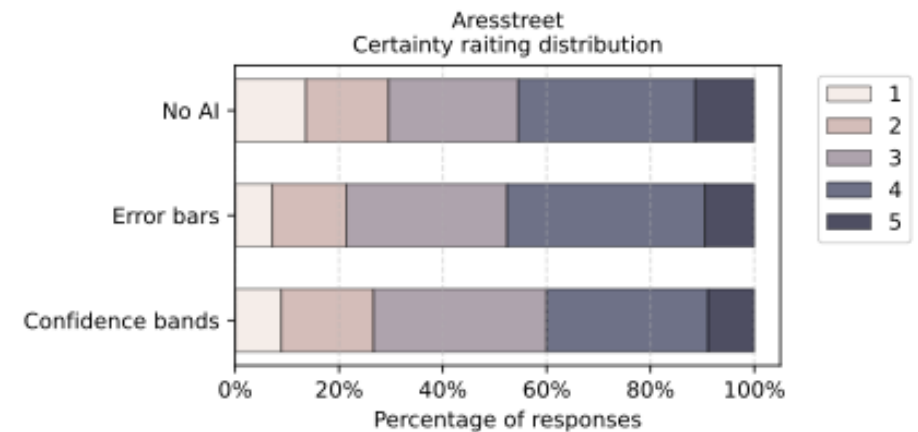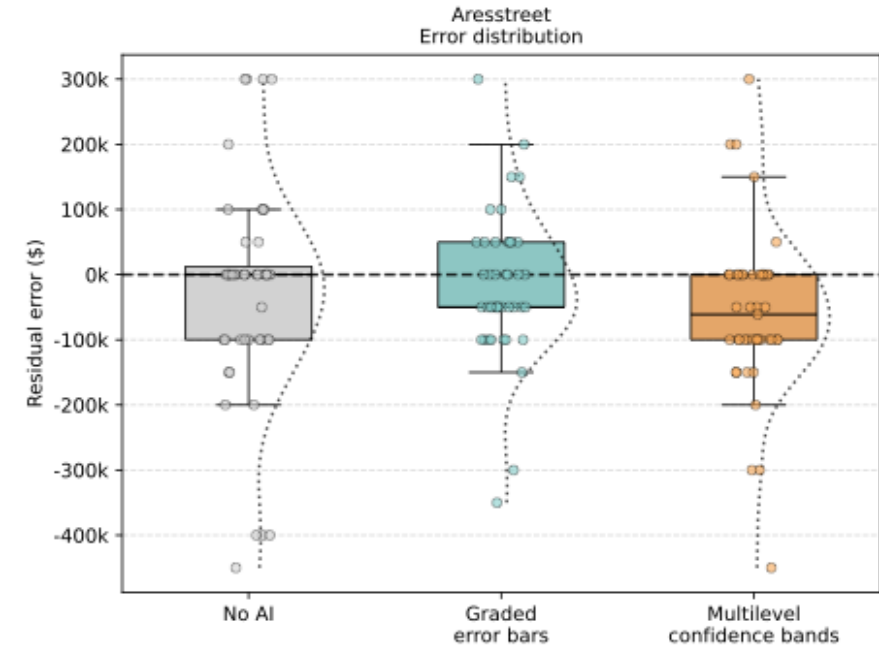# 6) Medium error – Aresstreet

# 6) Medium error – Aresstreet



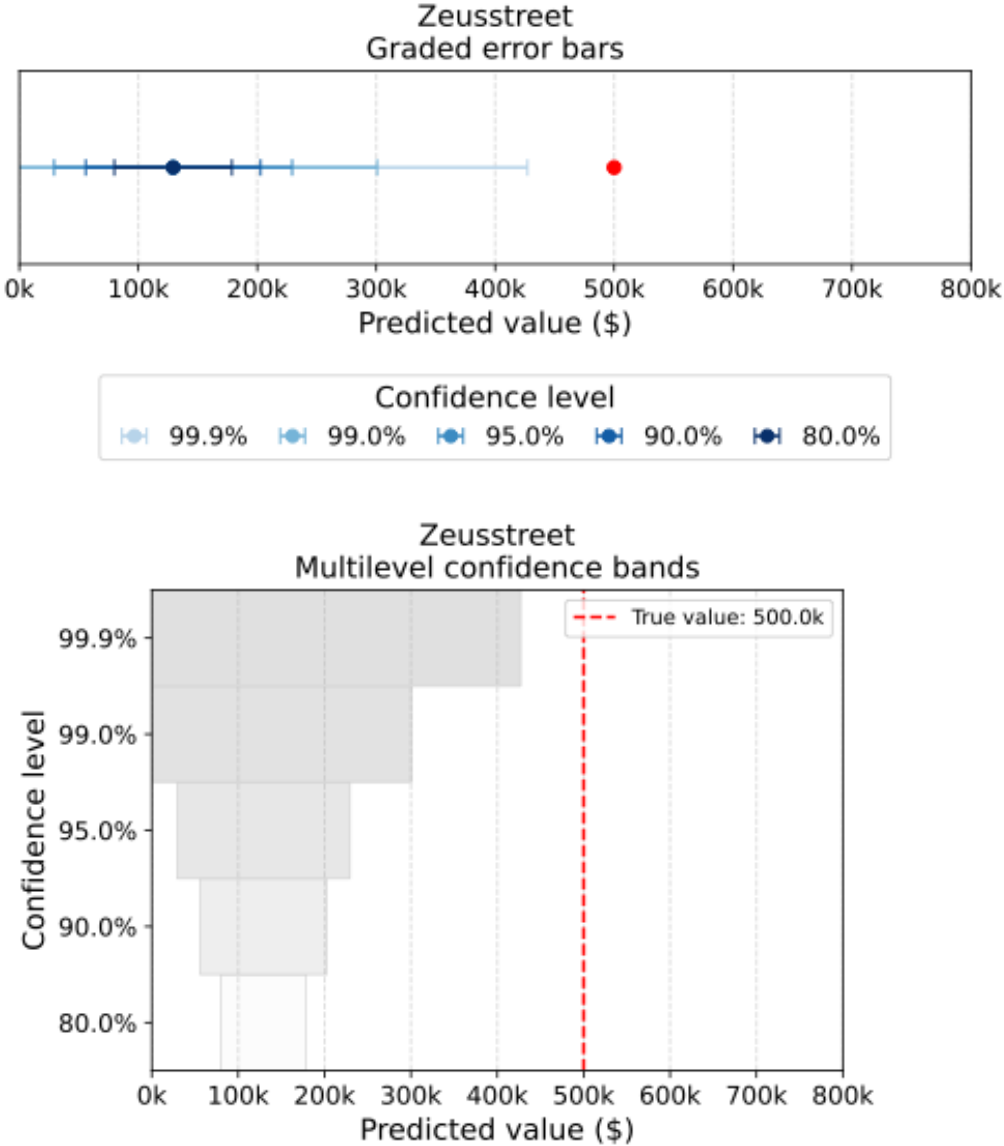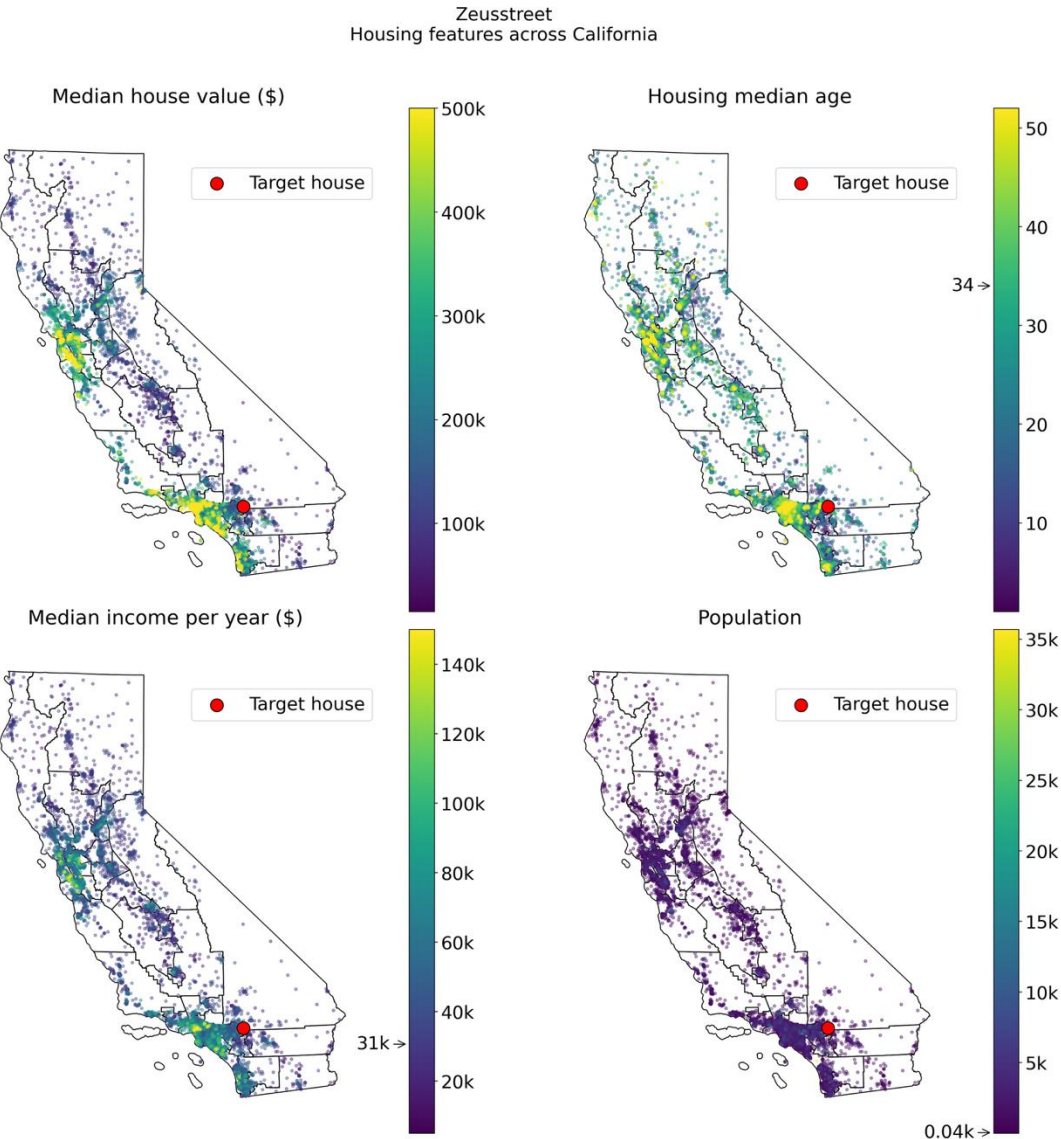Estimates closest to the true value are marked with an "x" above the bars.

# Task complexity: hard

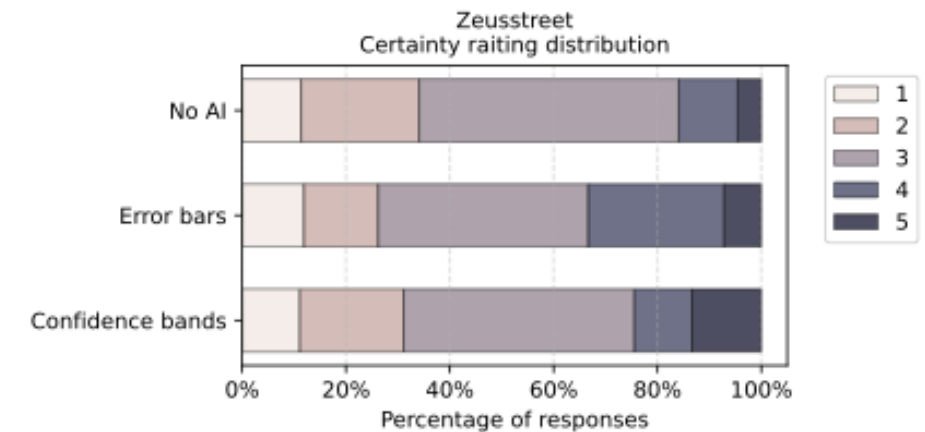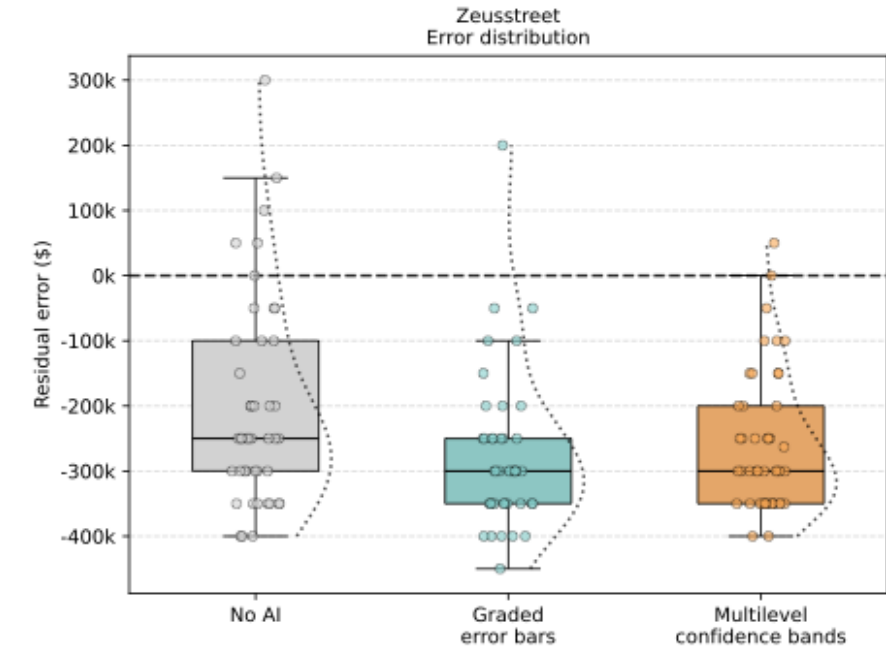- 🏠 7) Zeusstreet
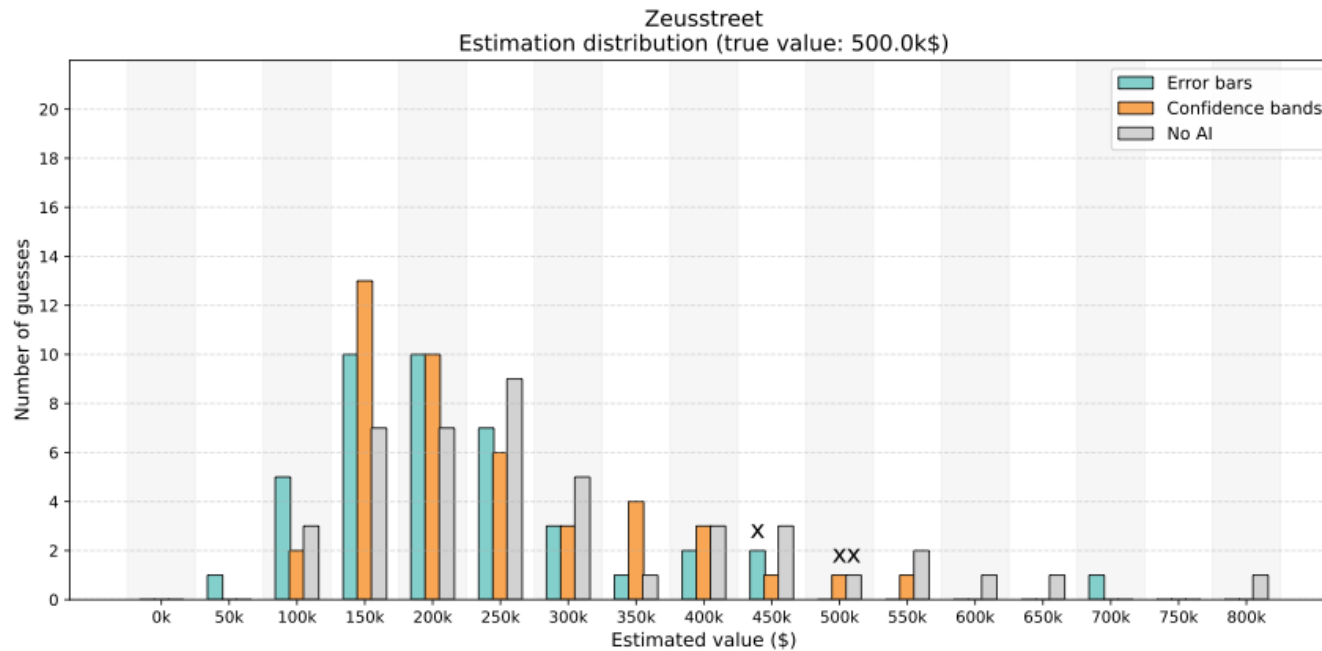- 🏠 8) Hephaestusstreet
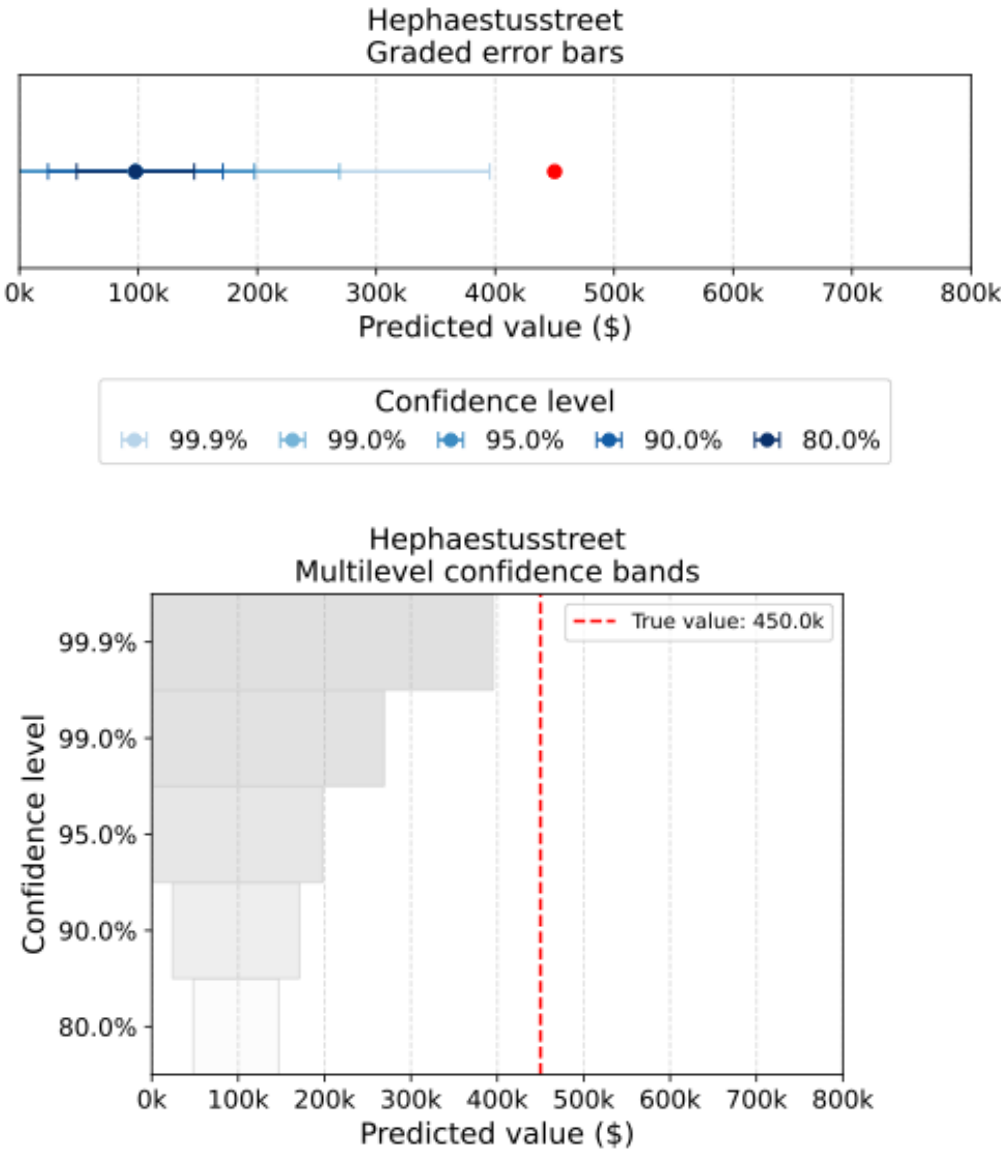- 🏠 9) Hestiastreet

# 7) High error – Zeusstreet



Estimates closest to the true value are marked with an "x" above the bars.
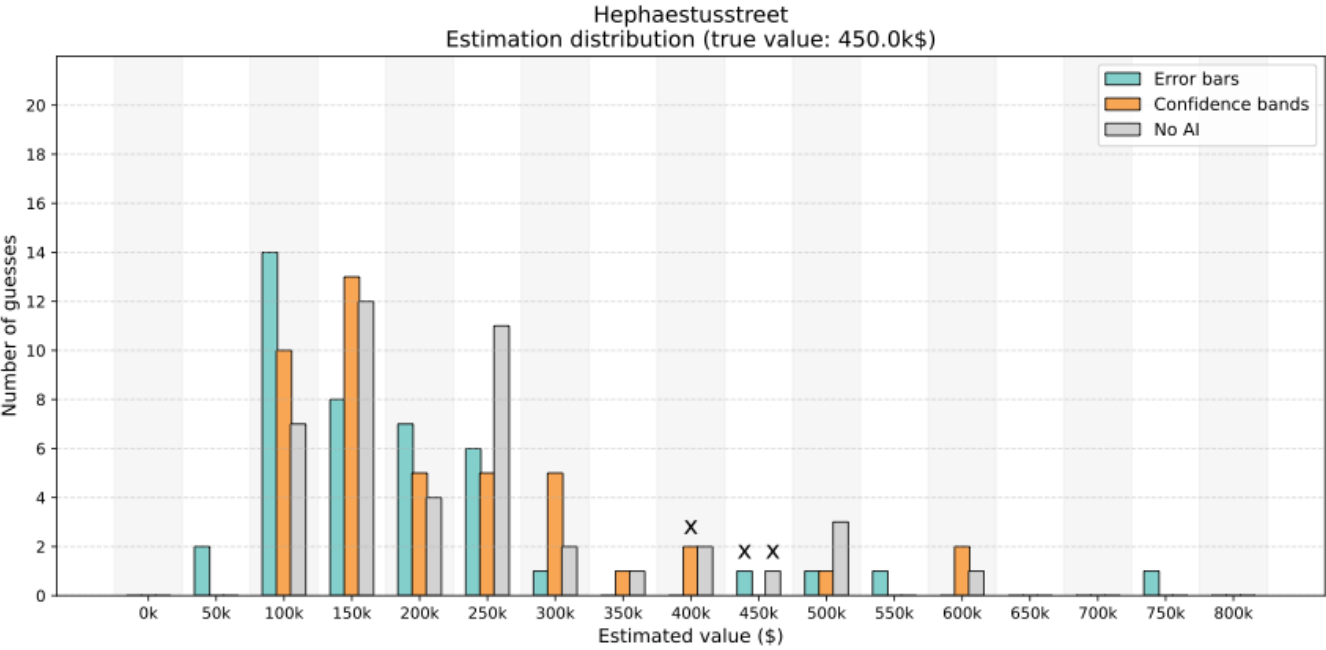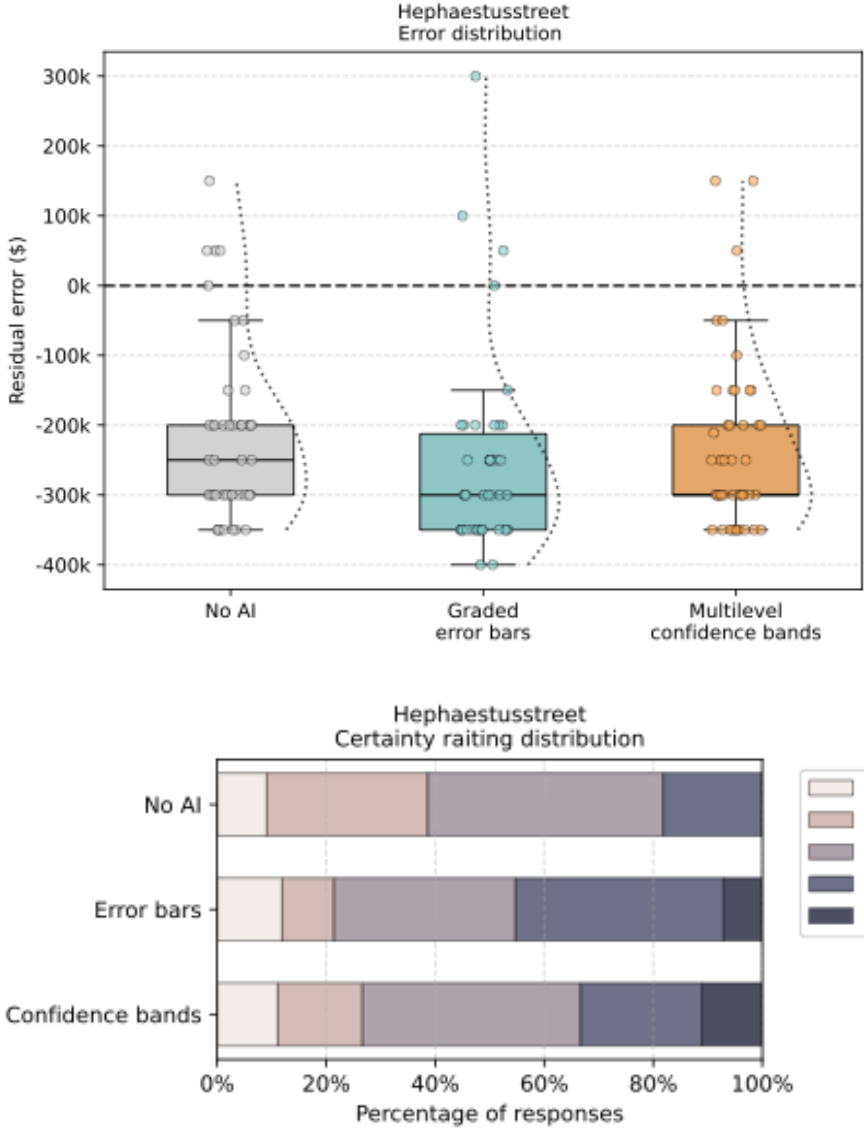
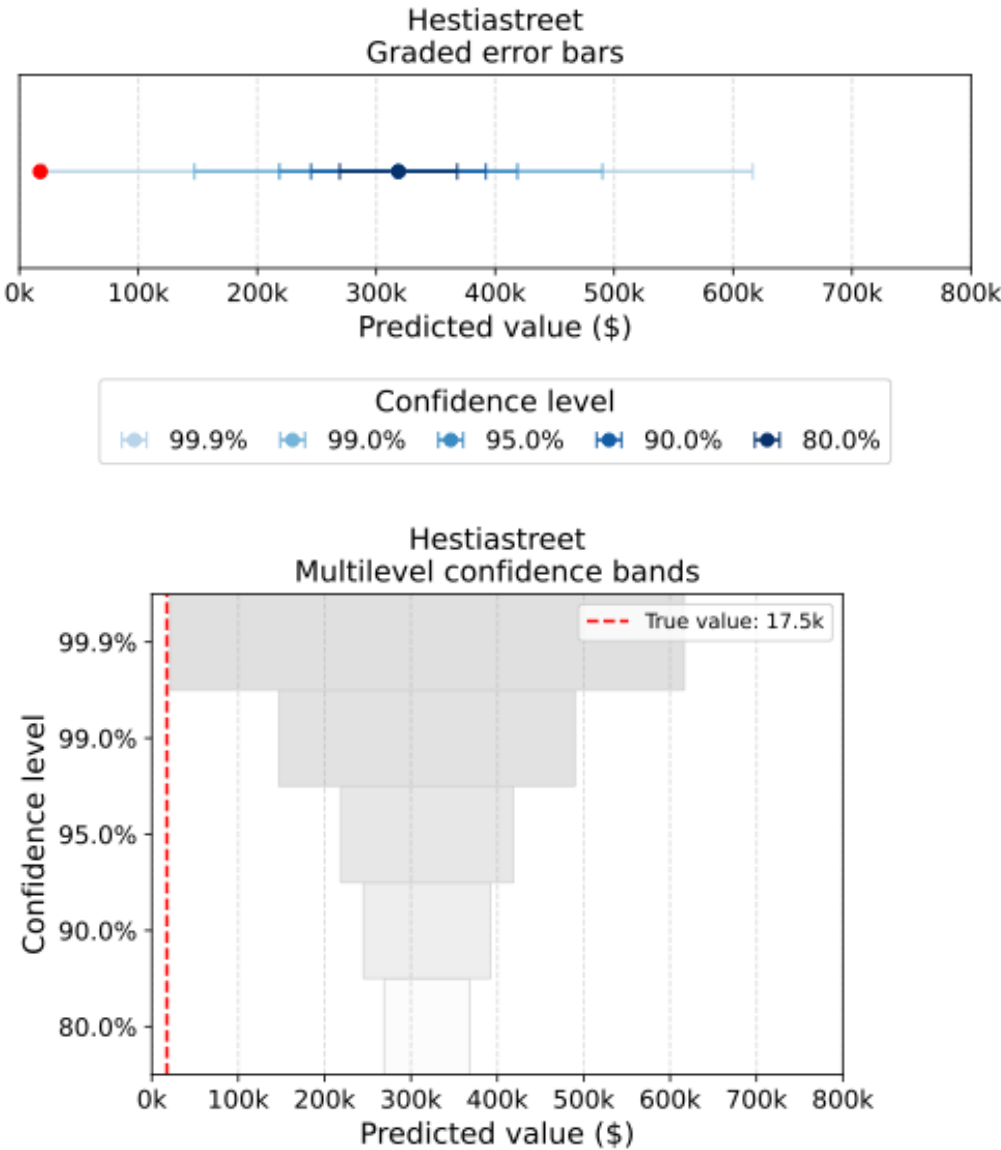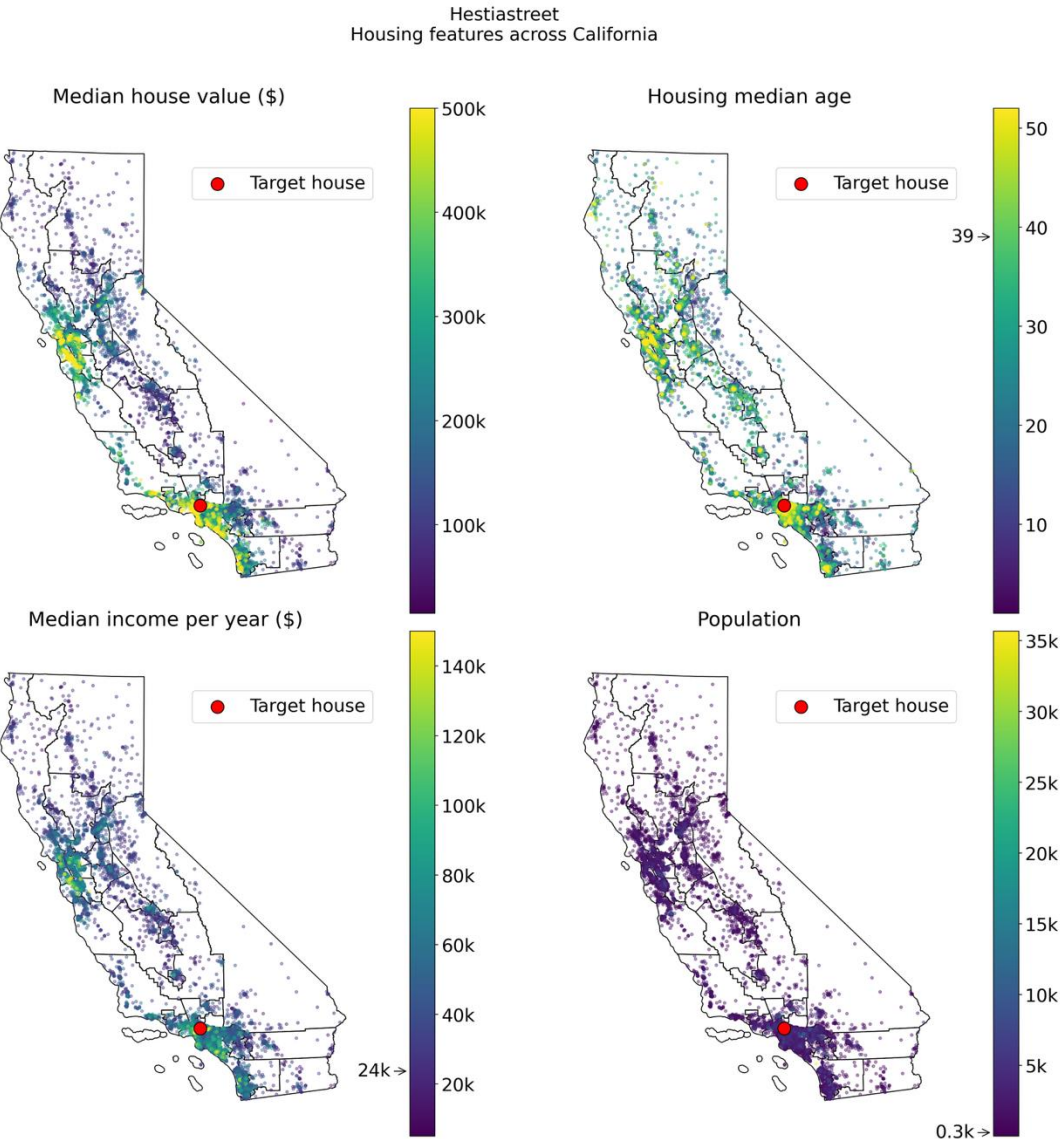# 8) High error – Hephaestusstreet

# 8) High error – Hephaestusstreet



Estimates closest to the true value are marked with an "x" above the bars.

# 9) High error – Hestiastreet

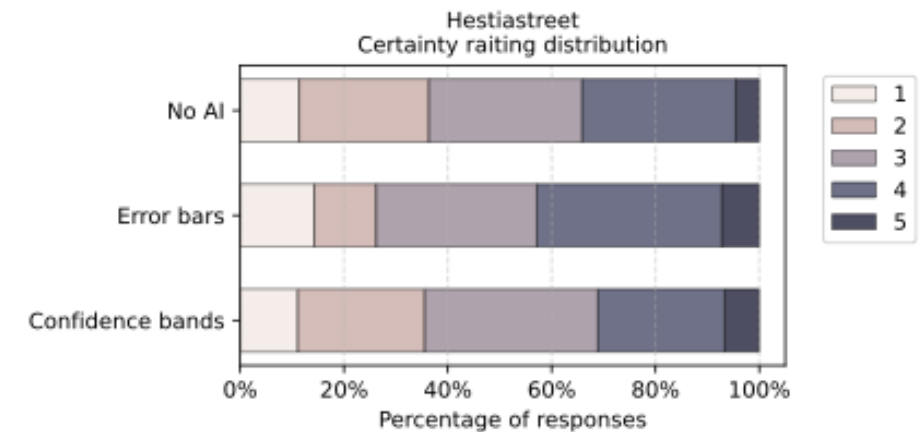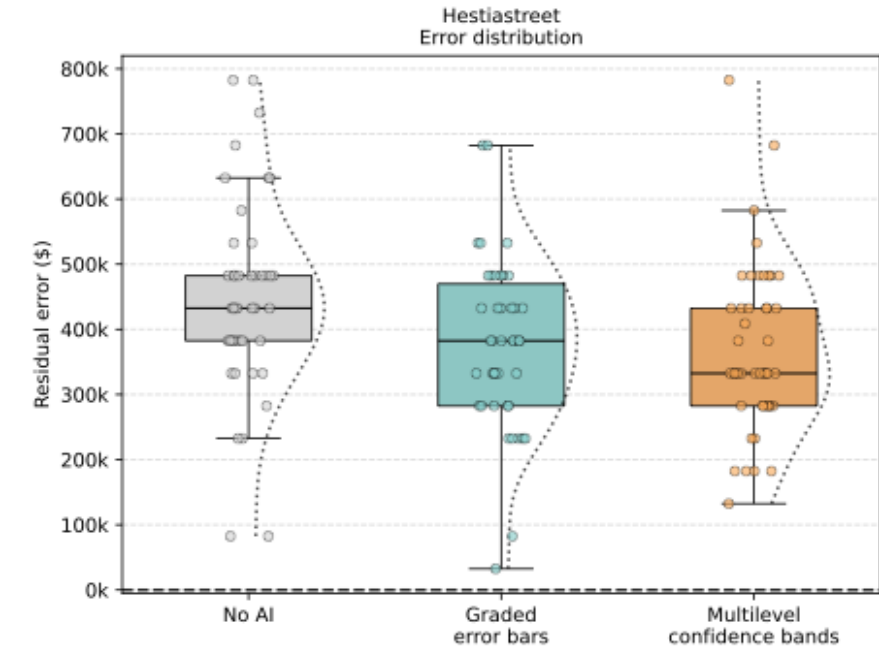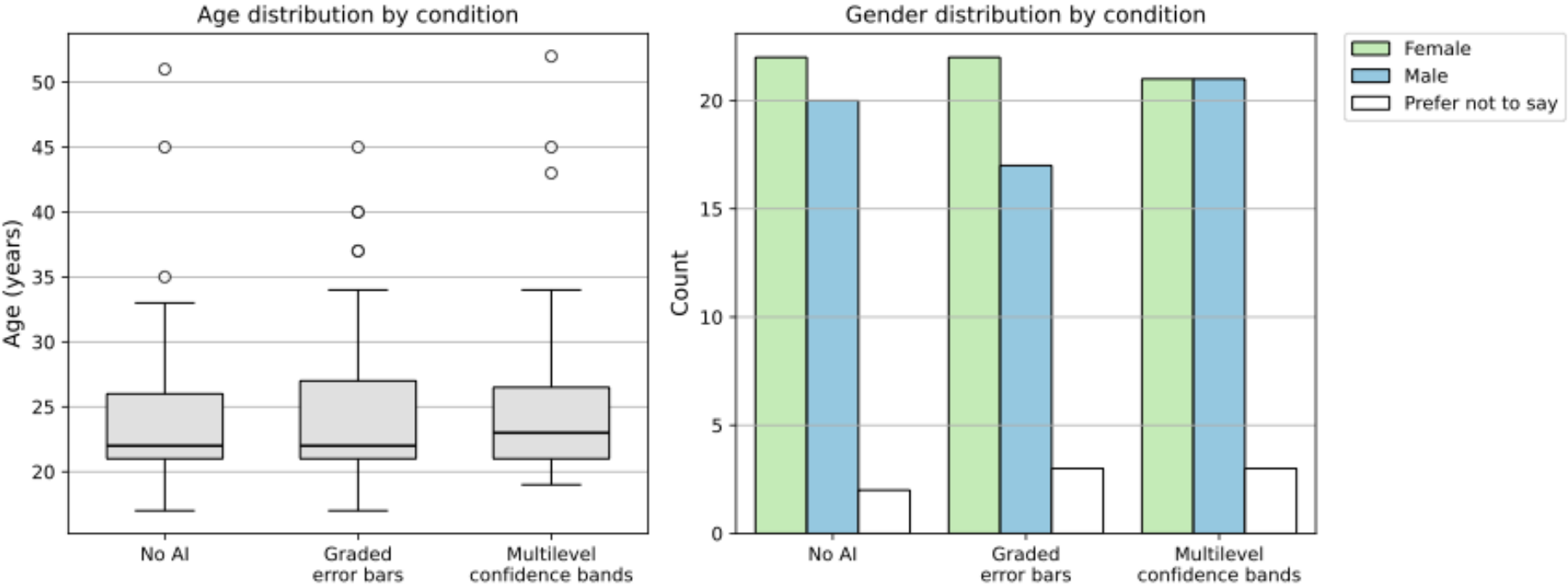

Estimates closest to the true value are marked with an "x" above the bars.
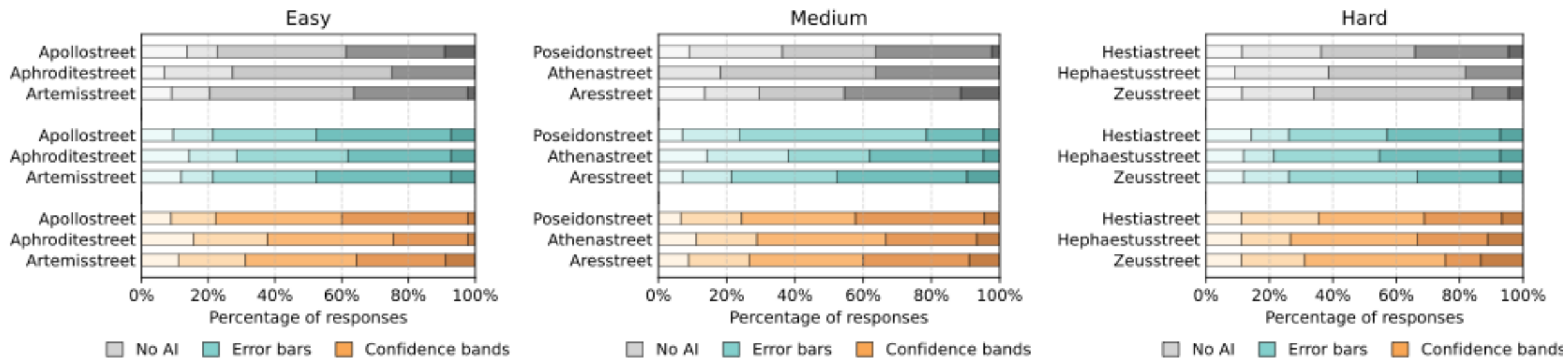
# 2. Additional charts

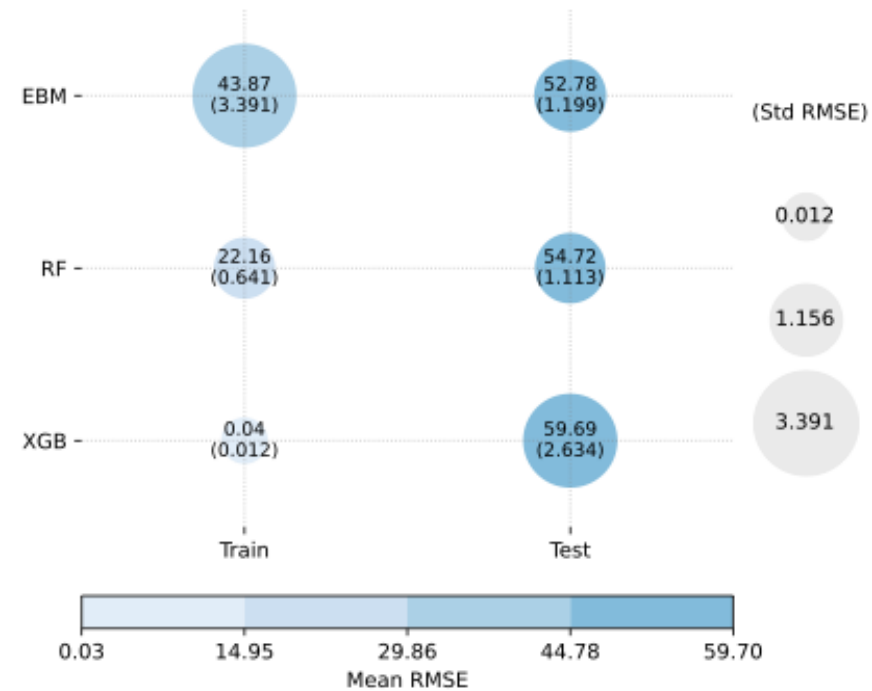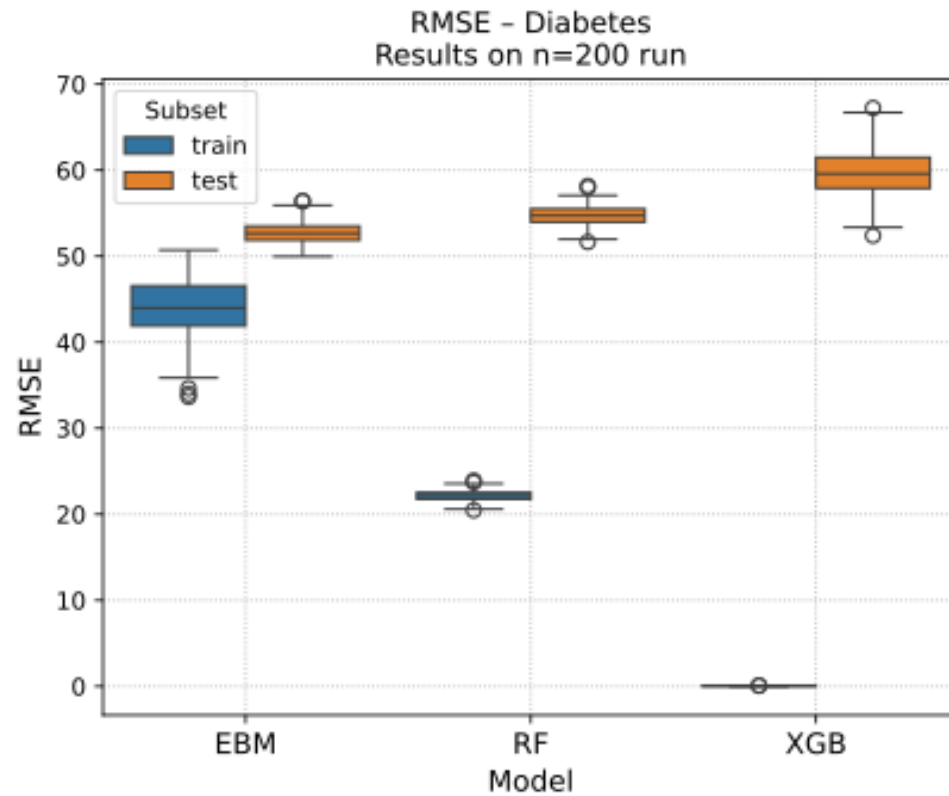# Survey Data: Demographic Statistics



Demographic characteristics

# Survey Results: Certainty Ratings by House, Visualization Condition and Complexity Level
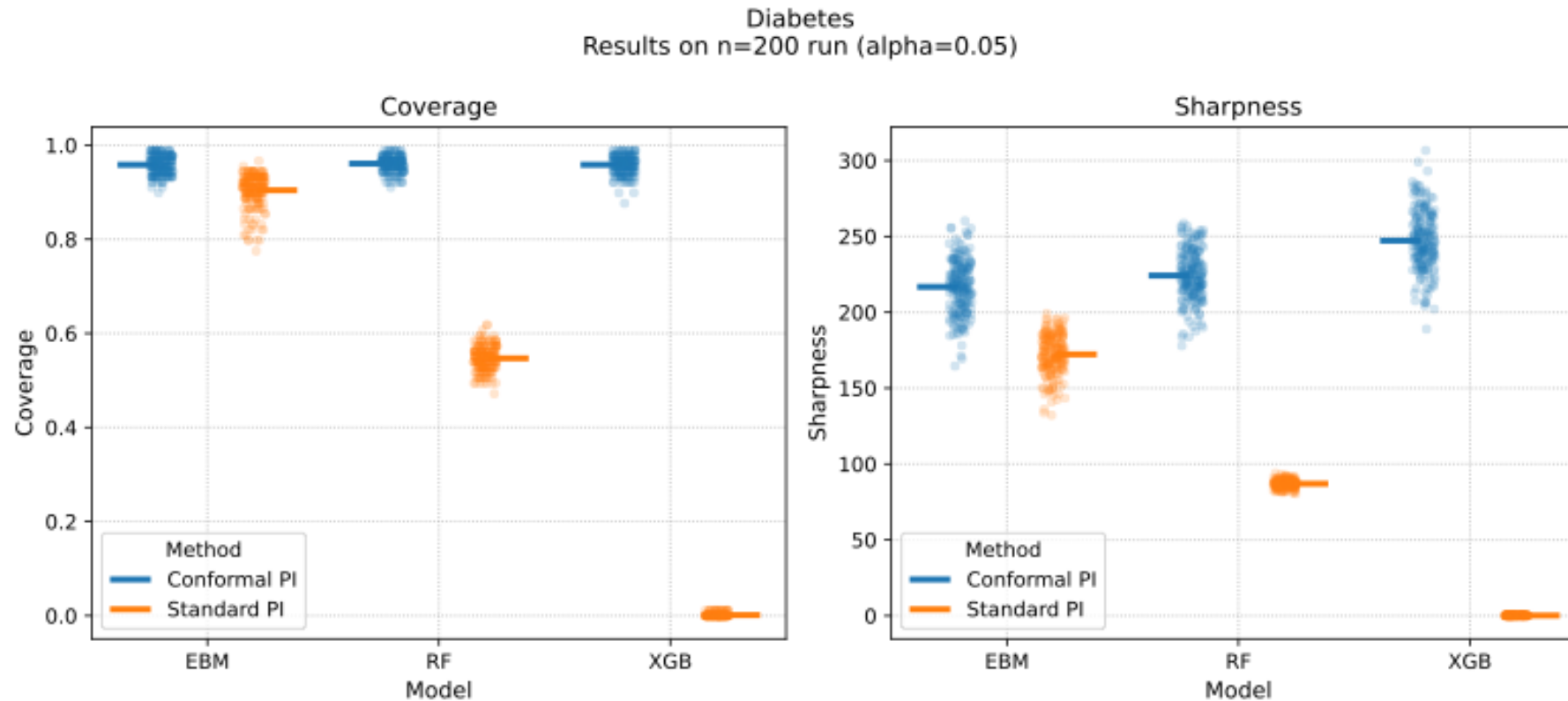
# Methods: Simulation Study – RMSE on *Diabetes*

The two visualizations display the same RMSE results across 200 Monte Carlo runs.
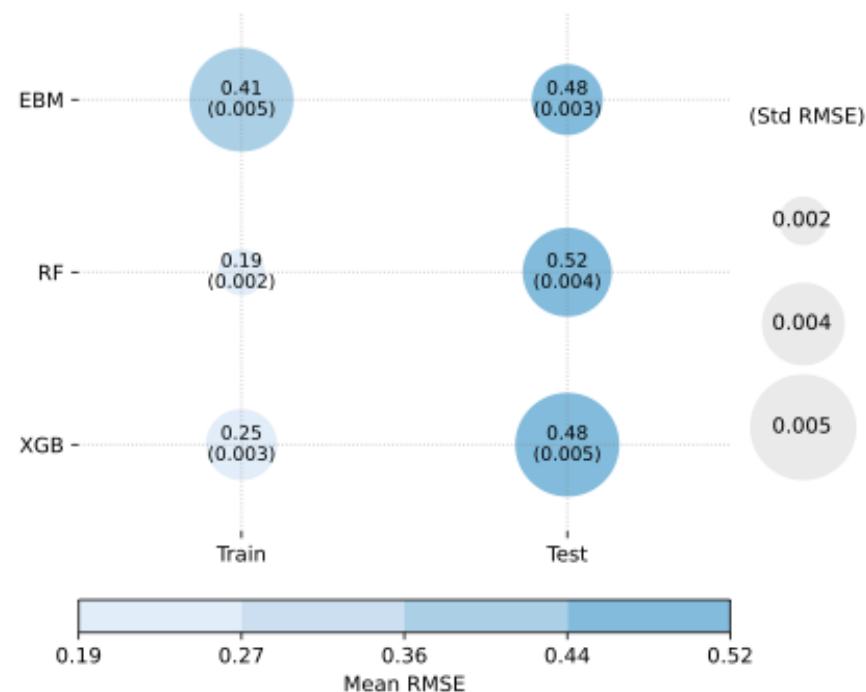
# Methods: Simulation Study – *Coverage* and *Sharpness* of Conformal vs Standard PIs on *Diabetes*



Diabetes
Results on n=200 run (alpha=0.05)

# Methods: Simulation Study – RMSE on *California Housing*

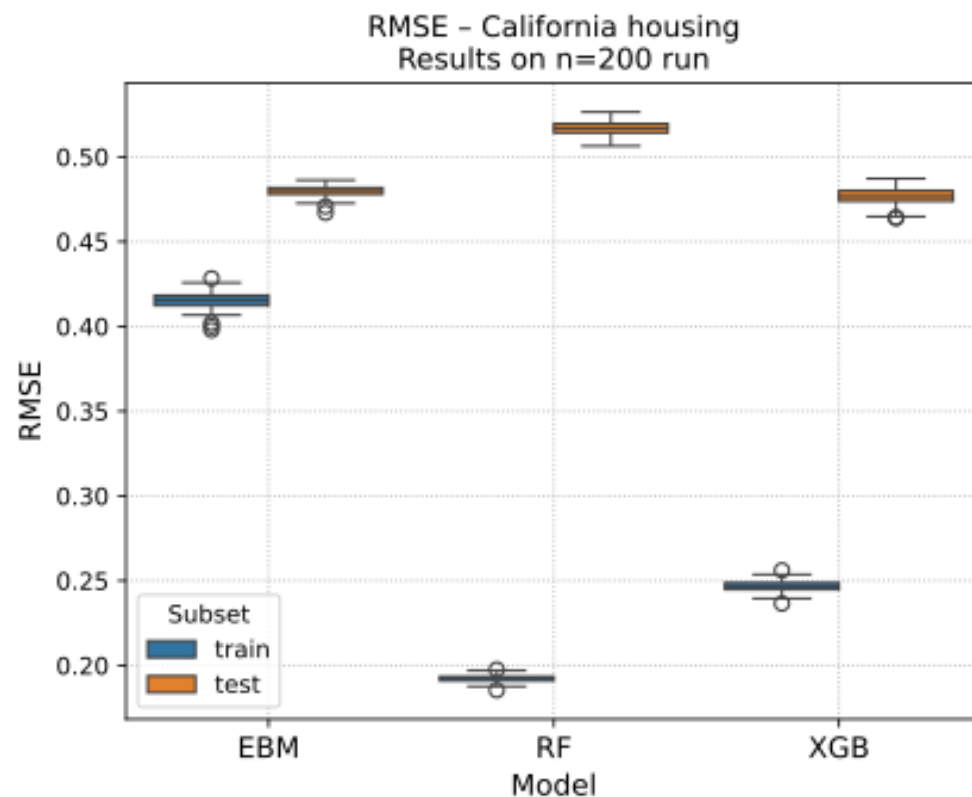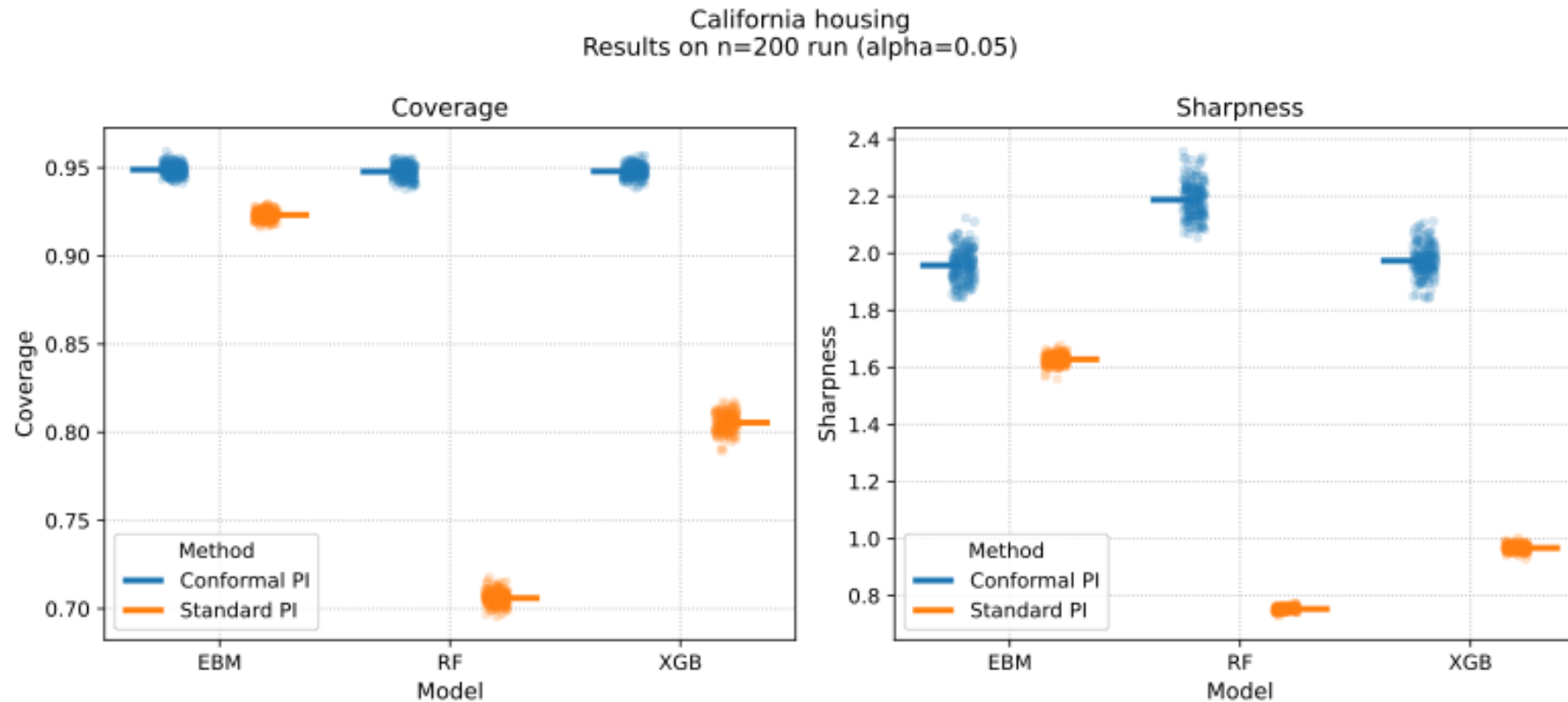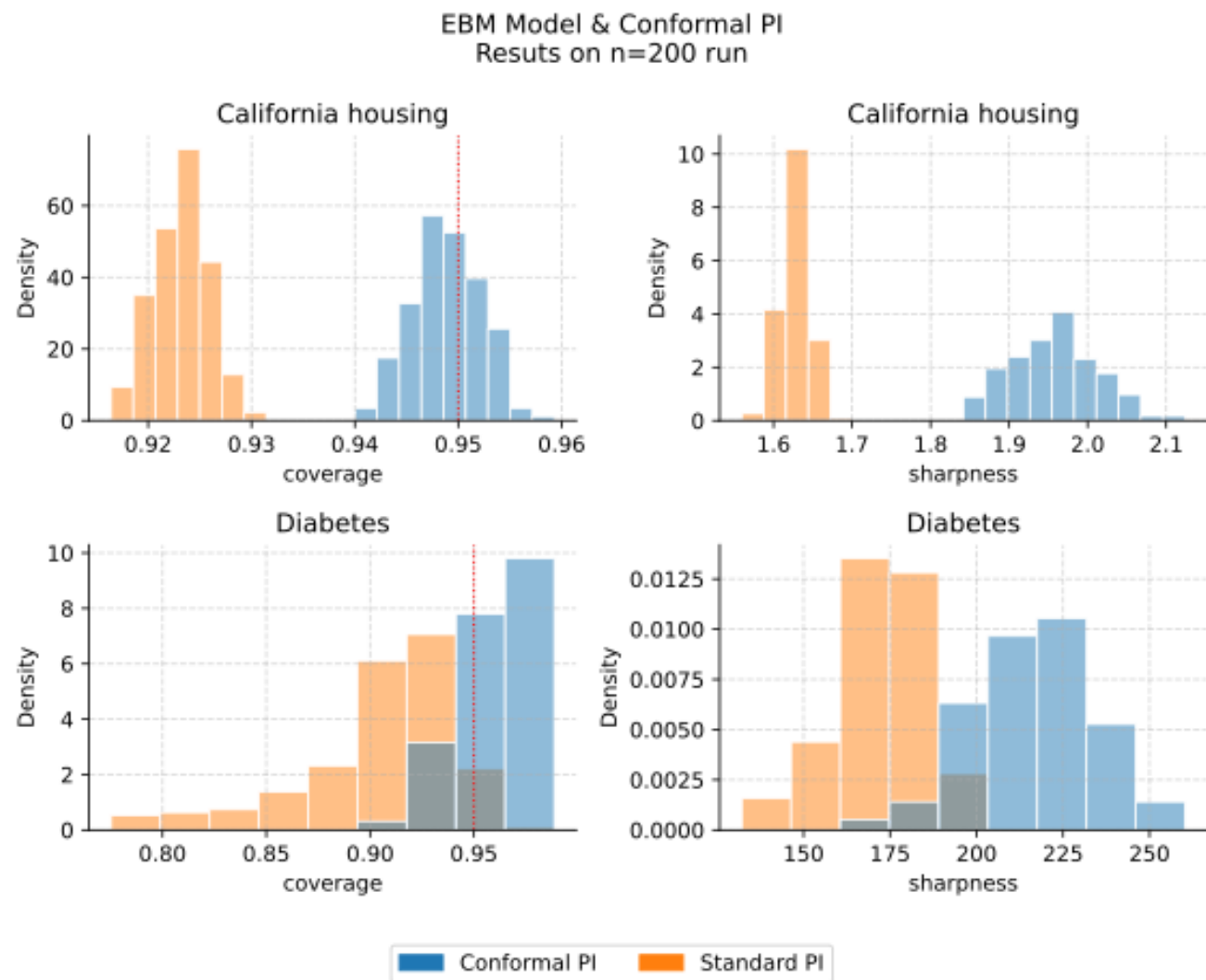The two visualizations display the same RMSE results across 200 Monte Carlo runs.

# Methods: Simulation Study – *Coverage* and *Sharpness* of Conformal vs Standard PIs on *California Housing*



California housing
Results on n=200 run (alpha=0.05)

# Methods: Simulation Study – Conformal PIs Statistics (on EBM)

# Coverage Analysis: Simulation Study

*Coverage* and *sharpness* evaluated over 200 Monte Carlo simulations across five α levels (Section 4.2).