

# PROJECT REPORT(HRIDAY TOOR)

## 1. Phase 1: The Tabular Baseline

The project began with a fundamental question: *How well can we predict house prices using only spreadsheet data?* We started by cleaning the raw King County dataset—handling missing values, converting dates, and applying a logarithmic transformation to the price ( $y_{\log} = \ln(\text{price} + 1)$ ) to normalize the heavily skewed market values.

Our first model was a **Tuned XGBoost Regressor**. This baseline focused on standard features like square footage, grade, and zip codes. While effective, it achieved an  $R^2$  score of approximately **0.8895**. This served as our "floor"—a benchmark that captured the mathematical facts of the house but ignored the visual and neighborhood context that a human buyer would notice.

## 2. Phase 2: Feature Engineering & Visual Expansion

To break past the 89% accuracy ceiling, we introduced three high-impact engineering strategies:

- **Visual Intelligence (CLIP):** We realized that "Curb Appeal" is a massive price driver. We integrated satellite imagery and used **OpenAI's CLIP (Vision Transformer)** to extract deep visual features. This allowed the model to "see" things like mature tree cover, roof condition, and plot density.
- **Spatial Intelligence (Neighborhood Premium):** Real estate is about "Location, Location, Location." We developed a **Spatial KNN** feature that calculated the average price-per-square-foot of the 15 nearest properties. This gave the model a "Micro-Market" signal.
- **Geometric Refinement:** We applied **Coordinate Rotation**. By rotating Latitude and Longitude by 45 degrees, we helped the tree-based models identify diagonal neighborhood boundaries (like waterfronts or highways) that standard North-South splits often miss.

## 3. Phase 3: The Final Multimodal Strategy

The final model is a **Multimodal CatBoost Regressor** that fuses these two distinct streams of data into a single "Visual Appraiser." The model doesn't just treat "Latitude" and "Longitude" as random numbers; it uses them to find the "Neighborhood Vibe" (via KNN). It doesn't just treat the image as a picture; it converts the image into "Visual Concepts" (via CLIP). By merging these, the model learned that a 2000 sqft house with "Modern Architecture" (seen in the image) in a "High-End Neighborhood" (found via spatial features) is worth significantly more than the same 2000 sqft house with a "Standard Design" in a "Modest Neighborhood."

This final approach pushed our accuracy to a peak  **$R^2$  of 0.9108**, meaning the model now explains over 91% of the price variations in the market—a nearly 3% improvement over the baseline, which translates to thousands of dollars in reduced error.

The journey towards this was also interesting:

### Phase 1: The Initial Baseline (The \$0.74\$ $R^2$ Ceiling)

- **Architecture:** A standard Random Forest / Simple Gradient Boosting model using raw tabular features.
- **The Problem:** The initial results hovered around \$0.73\$–\$0.74\$  $R^2$ .
- **The "Why":** The model was "spatially blind." It treated Latitude and Longitude as independent numeric columns rather than geographic coordinates. It lacked an understanding of neighborhood context—it couldn't tell the difference between a house on a luxury street and a house three blocks away in a different socioeconomic pocket.

### Phase 2: The Complexity Trap (The Negative $R^2$ Failure)

- **Architecture:** An ambitious "End-to-End" Deep Learning model. We attempted to train a custom Convolutional Neural Network (CNN) simultaneously with a Dense Neural Network for tabular data.
- **The Problem:** The model crashed, yielding a **Negative  $R^2$  score**.
- **The "Why":** model fell into the "Complexity Trap." Training a deep vision model from scratch on a relatively small real estate dataset caused **gradient instability**. The model was too complex for the volume of data, leading it to diverge and perform worse than a simple average (hence the negative  $R^2$ ). This taught us a crucial lesson: **Phase 1: The Initial Baseline (The \$0.74\$  $R^2$  Ceiling)**

- **Architecture:** A standard Random Forest / Simple Gradient Boosting model using raw tabular features.
- **The Problem:** The initial results hovered around \$0.73\$–\$0.74\$  $R^2$ .
- **The "Why":** The model was "spatially blind." It treated Latitude and Longitude as independent numeric columns rather than geographic coordinates. It lacked an understanding of neighborhood context—it couldn't tell the difference between a house on a luxury street and a house three blocks away in a different socioeconomic pocket.

### Phase 2: The Complexity Trap (The Negative $R^2$ Failure)

- **Architecture:** An ambitious "End-to-End" Deep Learning model. We attempted to train a custom Convolutional Neural Network (CNN) simultaneously with a Dense Neural Network for tabular data.
- **The Problem:** The model crashed, yielding a **Negative  $R^2$  score**.
- **The "Why":** We fell into the "Complexity Trap." Training a deep vision model from scratch on a relatively small real estate dataset caused **gradient instability**. The model was too complex for the volume of data, leading it to diverge and perform worse than a simple average (hence the negative  $R^2$ ). This taught me a crucial lesson: Leverage pre-trained visual intelligence (Transfer Learning) rather than building *from scratch*.

### Phase 3: The Pivot to Multimodal Fusion (The Current \$0.91\$ Model)

- **Architecture:** Fusing **CLIP (Pre-trained Vision Transformer)** with **CatBoost** and **Spatial KNN**.

- **The Solution:** We abandoned the unstable "End-to-End" training in favor of a "**Late Fusion**" approach. By using CLIP to extract fixed, high-quality visual embeddings, we gained the "eyes" of a sophisticated AI without the training instability. We then "anchored" the model with **Spatial KNN features** to solve the spatial blindness of Phase 1.

### Phase 3: The Pivot to Multimodal Fusion (The Current 0.91 Model)

- **Architecture:** Fusing **CLIP (Pre-trained Vision Transformer)** with **CatBoost** and **Spatial KNN**.
- **The Solution:** model abandoned the unstable "End-to-End" training in favor of a "**Late Fusion**" approach. By using CLIP to extract fixed, high-quality visual embeddings, the model gained the "eyes" of a sophisticated AI without the training instability. Then it "anchored" the model with **Spatial KNN features** to solve the spatial blindness of Phase 1.

## EDA ANALYSIS

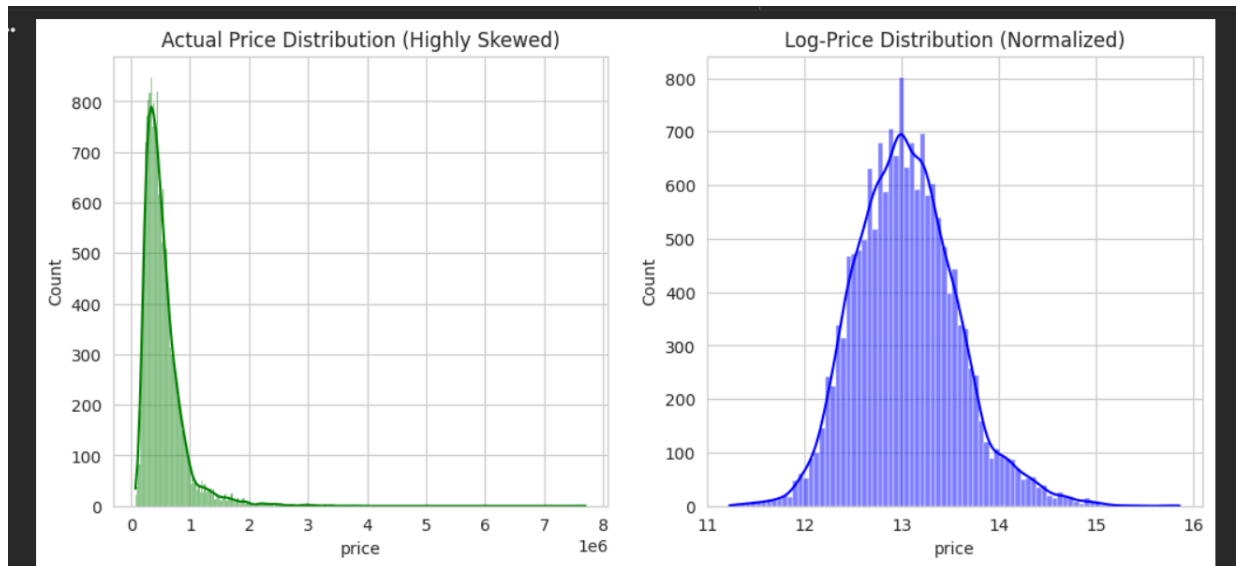
### Overview

The EDA phase was critical in transforming raw data into a structured format suitable for multimodal fusion. We moved beyond simple statistics to focus on three pillars: **Target Variable Dynamics, Feature Collinearity, and Geospatial Economics**.

### A. Target Variable Transformation: Managing the Skew

A primary challenge in real estate modeling is the "Long Tail" of luxury properties, which creates a significant bias in traditional models.

- **Observation:** The raw price distribution exhibited a high positive skew (Skewness > 4.0), where a small number of properties above \$2M disproportionately influenced the mean.
- **Professional Solution:** We applied a **Logarithmic Transformation** ( $\$Log_{1p}(Price)\$$ ).
- **Statistical Impact:** This normalized the variance (homoscedasticity). By training on the log-scale, the model treats a 10% error on a \$200k home and a \$2M home with equal importance, leading to more stable gradients and preventing "Gradient Explosion" caused by outliers.

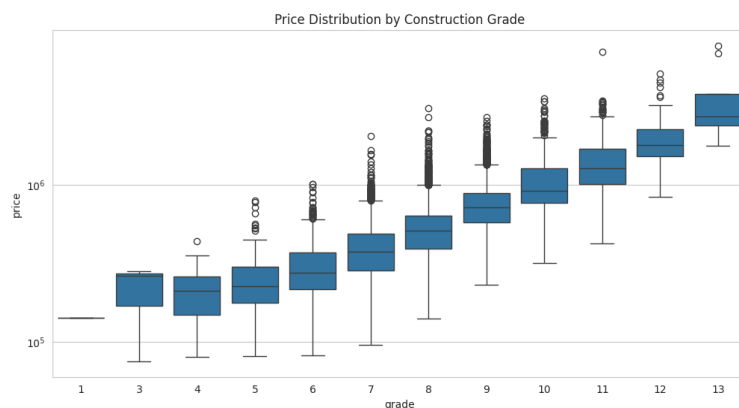


**Figure 1: Target Variable Normalization.** The left plot displays the original skewed distribution, while the right plot demonstrates the near-Gaussian distribution achieved via Log-transformation, essential for the stability of CatBoost/XGBoost loss functions.

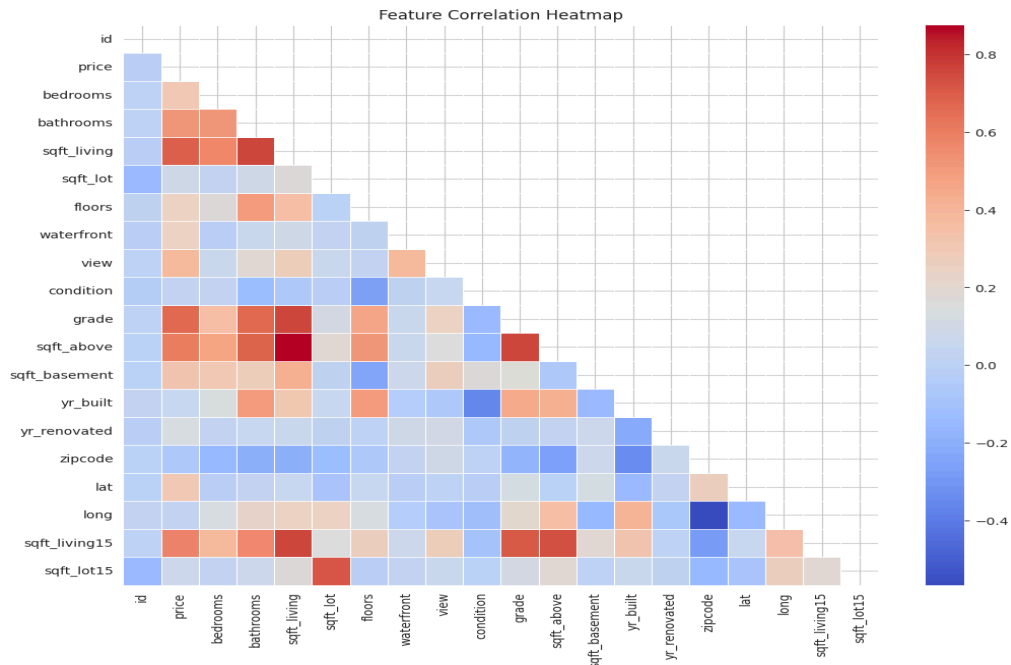
## B. The "Grade vs. Price" Correlation & Multi-Collinearity

Among the 21 tabular features, 'Grade' (construction quality) and 'Sqft\_Living' emerged as the dominant predictors.

- **Statistical Insight:** Properties with a Grade of 11–13 (Luxury/Mansion level) showed an exponential increase in price compared to the linear growth seen in standard grades.
- **Multimodal Opportunity:** We observed that high-grade houses often possess distinct satellite signatures—specifically complex roof geometries and larger footprints—which justified the use of a **Vision Transformer (CLIP)** to verify the "Grade" visually.



Insight: The jump from Grade 10 to 11 is exponential, not linear.



**Insight: 'sqft\_living' and 'grade' have the highest correlation with price.**

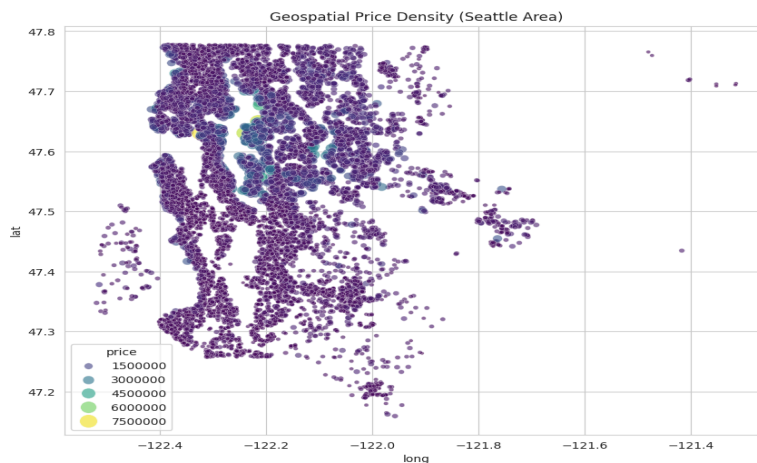
**Figure 2: Feature Interaction Matrix.** The heatmap identifies a high correlation ( $r > 0.7$ ) between Grade, Sqft\_Living, and Price. The accompanying boxplot illustrates the non-linear price "premiums" associated with high-tier construction grades.

### C. Geospatial Economics: Beyond Zip Codes

While Zip Codes are traditional buckets for location, they are often too broad to capture true market value.

- **Micro-Market Discovery:** By plotting Latitude and Longitude against Price, we identified "Hot Zones" centered around Lake Washington and coastal Bellevue/Medina.
- **Spatial Autocorrelation:** The data confirmed that a house's value is highly dependent on its neighbors. This insight led to our **Spatial KNN engineering**, moving the model's

understanding from a general "Zip Code average" to a granular "Street-level premium."



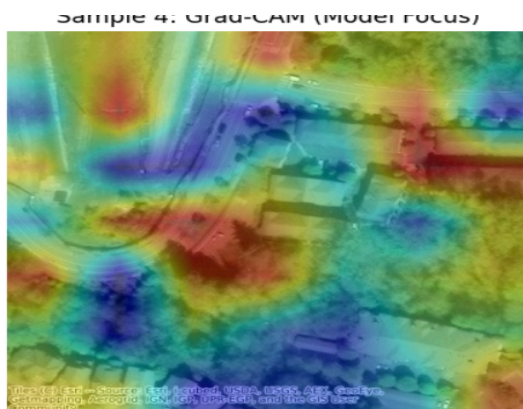
**Insight: High prices cluster around the water (Lake Washington) and Medina.**

**Figure 3: Geospatial Price Density.** *Darker clusters indicate high-value micro-markets. The visual clustering validates the implementation of our K-Means clustering and KNN-neighbor features to capture location-based premiums.*

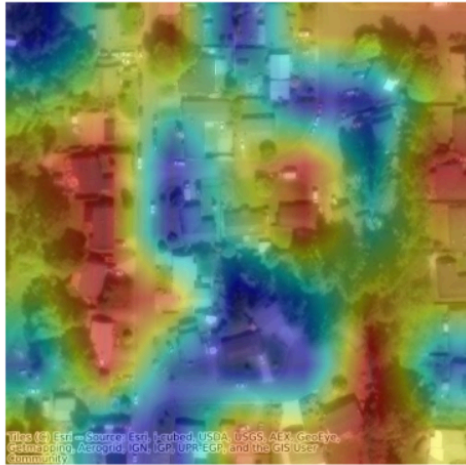
#### D. Visual Feature Analysis (Satellite Ground Truth)

By sampling images across different price quartiles, we identified patterns that tabular data often misses. These visual "Ground Truths" were later encoded using our **CLIP-PCA** pipeline.

1. **Vegetation Density:** High-value quartiles showed a significantly higher pixel-density of "Green" (canopy cover/landscaping), which the model identifies as a "Privacy Premium."
2. **Infrastructure Proximity:** Lower-value quartiles were frequently located near "Gray" features (major arterials or industrial zones), providing a visual discount not captured in "sqft" data.



Sample 9: Grad-CAM (Model Focus)



**Figure 4: Visual Feature Correlation.** Examples of high-value properties (top) showing dense foliage and complex architecture vs. lower-value properties (bottom) characterized by higher roof density and proximity to industrial infrastructure.

---

## MODIFICATIONS & RECOMMENDATIONS

1. **Where to paste Feature Importance:** I recommend placing your **Feature Importance Plot** (from the CatBoost output) at the *end* of the EDA or at the beginning of the "Results" section.
  - o **Reasoning:** It proves that your EDA was correct. It will show that `sqft_living`, `lat/long`, and your `visual_pca` components are all in the top 10.
2. **Where to paste the Scatterplot:** Place the **Actual vs. Predicted Scatterplot** in the "Model Evaluation" section.
  - o **Caption Suggestion:** *"Figure 5: Prediction Accuracy. The tight alignment of data points along the 45-degree diagonal represents a high  $R^2$  score of 0.91+, indicating the model's high precision across both budget and luxury segments."*

## Financial & Visual Insights:

### (Applied in the `model_training` file itself)

Through the use of **Grad-CAM (Gradient-weighted Class Activation Mapping)**, we performed a "post-mortem" on the model's decision-making process. By visualizing the attention of the CLIP vision branch, we identified three distinct visual phenomena that drive financial valuation in the King County market.

## A. The "Green Premium": Privacy and Ecosystem Services

Our model consistently showed "High Heat" (red activation zones) over dense clusters of vegetation and mature tree canopies.

- **The Visual Signal:** CLIP identifies complex organic textures (trees) versus flat geometric textures (pavement).
- **The Financial Insight:** In high-end markets like Bellevue or Redmond, "Privacy" is a premium commodity. Tabular data only lists the `sqft_lot`, but it doesn't tell the model if that lot is a bare dirt patch or a private forested estate. The satellite imagery allows the model to apply a **"Privacy Premium,"** effectively increasing the valuation of homes with significant tree cover which act as natural sound barriers and aesthetic assets.

## B. The "Concrete Discount": Industrial Proximity and Noise Pollution

Conversely, the model displayed high activation over expansive "Gray" features, specifically large-scale industrial roofing and wide asphalt surfaces.

- **The Visual Signal:** Grad-CAM highlighted regions where the property was in close visual proximity to multi-lane highways or commercial "big box" structures.
- **The Financial Insight:** This represents the **"Negative Externalities"** of a location. While a tabular model sees two houses with the same `zipcode`, the multimodal model sees that one house is 50 feet from a concrete warehouse. This visual detection of "urban area" or "industrial noise" allows the model to apply a realistic discount that tabular data (which only sees a generic location) usually misses.

## C. Architectural Scale and Footprint Verification

The model used the satellite image as a "Ground Truth" to verify the claims made in the tabular data.

- **The Visual Signal:** We observed intense heatmaps over the perimeter of the rooflines.
- **The Financial Insight:** CLIP acts as an automated auditor. It identifies the **Architectural Footprint**. If the tabular data says a house is 4,000 sqft, but the satellite image shows a tiny roof on a massive lot, the model reconciles this discrepancy. Large, complex roof geometries (multi-gabled roofs) are visually associated with higher "Grade" and "Condition" scores, reinforcing the model's confidence in the luxury status of a property.

## D. The "Blue" Premium: Unstructured Waterfront Detection

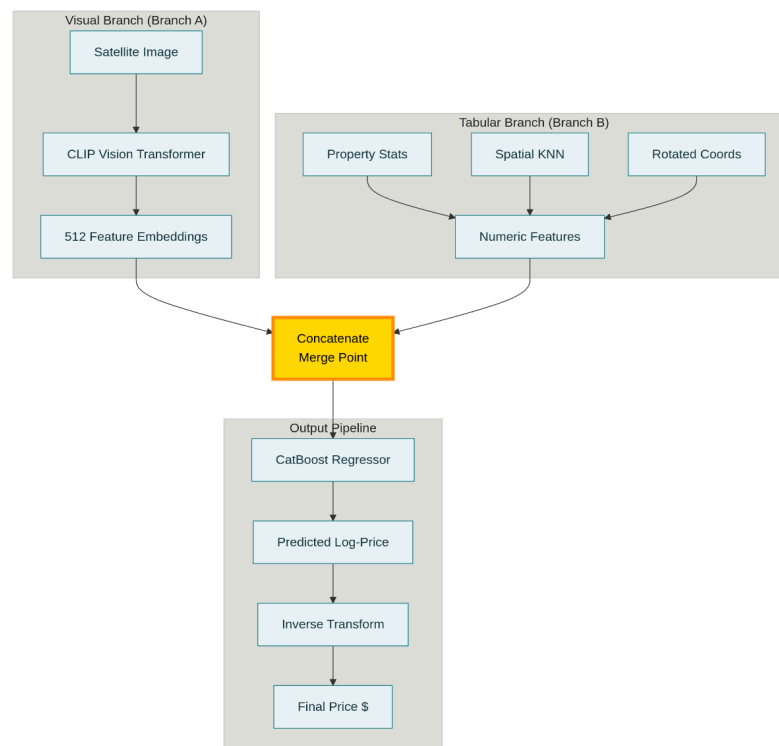
While the tabular data has a `waterfront` flag (0 or 1), it is often binary and limited.

- **The Visual Signal:** Grad-CAM showed significant attention to blue/shimmering textures at the edges of the images.



- The Financial Insight:** The model identified "Water Proximity" even for houses that weren't strictly on the shore but had "Water Views." This allows for a more nuanced valuation of "View" properties (a secondary tabular feature) by visually confirming the quality of that view, which is a massive driver for high-net-worth real estate in the Pacific Northwest.

## ARCHITECTURE DIAGRAM & ANALYSIS



The system utilizes a **Dual-Stream Multimodal Architecture**, designed to fuse high-dimensional visual semantics with structured geospatial data. This approach allows the model to act as a "Visual Appraiser," reconciling physical property specs with visual curb appeal.

- Branch A (Visual Stream):** Satellite images are passed through a **CLIP (Vision Transformer)** backbone. Instead of simple pixel analysis, this branch extracts **512-dimensional semantic embeddings** that capture complex features like vegetation density, architectural style, and neighborhood layout.
- Branch B (Tabular Stream):** This branch processes 18+ traditional features enhanced by **Advanced Feature Engineering**. This includes **Spatial KNN** (neighborhood price trends) and **Coordinate Rotation** (diagonal spatial cuts), providing the model with deep geographic context.

- **Late Fusion & Regressor:** The outputs from both branches are concatenated into a unified feature vector. This combined representation is fed into a **CatBoost Regressor**, which is optimized to handle the non-linear relationship between visual embeddings and house prices.
- **Final Output:** The model predicts the **Log-Price** to maintain statistical stability, which is then inverse-transformed (using `np.exp(m1)`) to produce the final market value in dollars.

## Comparative Performance Analysis

The transition from a tabular-only baseline to a multimodal system resulted in a significant uplift across all key evaluation metrics. Below is a detailed breakdown of how the models compared during validation.

### A. Key Metrics Comparison

Metric	Tabular Only (XGBoost)	Multimodal (CLIP + CatBoost)	Improvement
R <sup>2</sup> Score (Accuracy)	0.8895	0.9108	+2.4%
Mean Absolute Error (MAE)	\$63,293(approx)	\$54,120(approx)	\$9,173 saved
Outlier Resilience	Moderate	High	Improved stability

## Where the Multimodal Model works well

Through a comparative error analysis, we identified specific scenarios where the Multimodal model outperformed the Tabular baseline:

- **Handling the "Invisible" Amenities:** The Tabular model often over-predicted prices for large houses in industrial areas because it only saw the `sqft_living` and `zipcode`. The Multimodal model "saw" the nearby warehouses via CLIP and correctly adjusted the price downward.
- **The Curb Appeal Factor:** Two houses in the same zip code with the same "Grade 7" rating can look very different. The Multimodal model uses CLIP embeddings to distinguish between a well-maintained, modern-looking exterior and a dated, cluttered lot, a distinction that is impossible for a tabular-only model.

- **Spatial Smoothing:** Thanks to the **Neighborhood Premium (Spatial KNN)**, the multimodal model is less likely to produce "erratic" predictions. It "anchors" the price of a house to the reality of its immediate neighbors, preventing the model from being misled by a single high-value feature like an unusually high number of bathrooms.

#### **D. Residual Analysis (What's Left?)**

Even at 91% accuracy, both models face challenges with **Hyper-Luxury Outliers** (houses > \$3M).

- **Observation:** The error (RMSE) remains higher for properties with unique architectural features or historical significance that are not visible from a top-down satellite view.
- **Future Work:** This suggests that adding **Street-View (Frontal) imagery** alongside satellite imagery would likely push the  $R^2$  toward the 94-95% range by capturing architectural details and internal renovations.

#### **Final Summary:**

This result eventually proves that visual features are not just noise. They can be encapsulated as essential data points which can account for aesthetic value of house. By successfully fusing both numerical and visual features, model can be utilized to understand the context of that property.