

Satellite Imagery-Based Property Valuation

1. Overview: Approach and Modelling Strategy

Problem Statement:

Traditional real estate valuation models rely heavily on structured (tabular) attributes such as location, size, and amenities. However, they fail to capture visual and environmental cues(e.g., greenery, urban density) that strongly influence property prices.

This project proposes a multimodal learning approach that combines tabular property data with satellite imagery to improve price prediction performance.

Proposed Approach:

The overall strategy involves:

- * Using tabular features to capture structured information.
- * Using satellite images to extract neighbourhood-level visual features.
- * Learning joint representations by fusing outputs from both models.
- * Comparing performance against a tabular-only baseline.

Modelling Pipeline:

1. Data collection and preprocessing
2. Exploratory Data Analysis (EDA)
3. Tabular-only regression model
4. CNN-based image feature extraction
5. Multimodal fusion (Tabular + Image features)
6. Model evaluation and comparison

2. Exploratory Data Analysis (EDA)

This section explores the statistical, spatial, and relational characteristics of the dataset using structured visualisations. The objective is to understand how property prices vary with structural attributes, location, and overall housing quality, and to motivate the modelling choices used later.

Distribution of Property Prices:

The histogram of property prices reveals a highly right-skewed distribution, indicating that while most properties are concentrated in a lower to mid-price range, a small number of properties command extremely high prices.

Key observations:

- * Majority of houses fall well below the maximum observed price.
- * Presence of long tail suggests strong price heterogeneity.
- * Median price (marked using a dashed vertical line) lies significantly left of the maximum values, confirming skewness~



Log-Transformed Price Distribution:

To address skewness, a \log_{10} transformation of price was visualised.

Observations after transformation:

- * Distribution becomes more symmetric and bell-shaped.
- * Extreme outliers are compressed into a manageable range.
- * Improves numerical stability and learning efficiency for models.

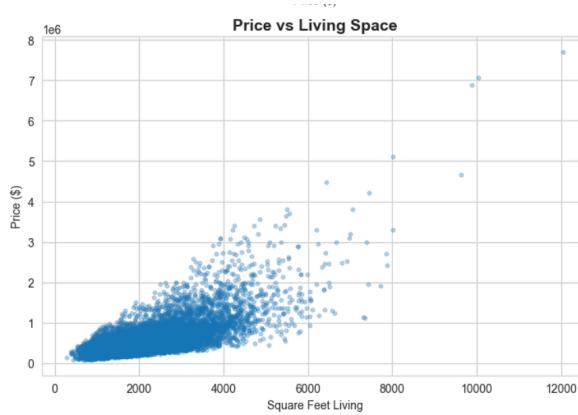


Relationship Between Living Area and Price:

A scatter plot was used to analyse the relationship between living area (sqft_living) and property price.

Key insights:

- * Strong positive relationship: prices generally increase with living space.
- * Significant vertical spread indicates that size alone does not fully determine price.
- * For similar square footage, prices vary widely — suggesting influence of other factors such as location, view, or neighbourhood quality.

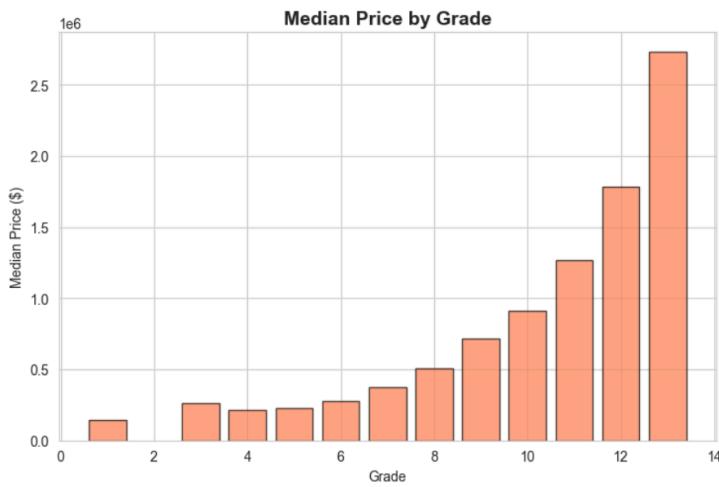


Impact of House Grade on Price:

House grade represents an overall assessment of construction quality and design. Median price was computed for each grade level.

Observations:

- * Median price increases monotonically with grade.
- * Higher grades correspond to disproportionately higher prices.
- * Indicates that grade is one of the strongest predictors of property value.

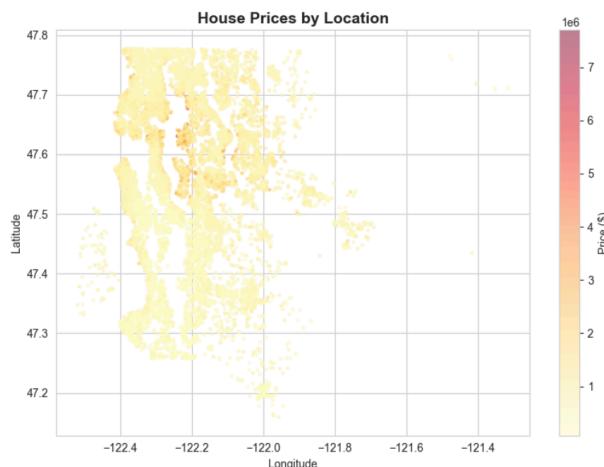


Geospatial Distribution of Property Prices:

A geographic scatter plot was created using longitude and latitude, with colour intensity representing price.

Key findings:

- * Clear spatial clustering of high-priced properties.
- * Certain geographic regions consistently exhibit higher prices.
- * Prices are not uniformly distributed, confirming strong location dependency.

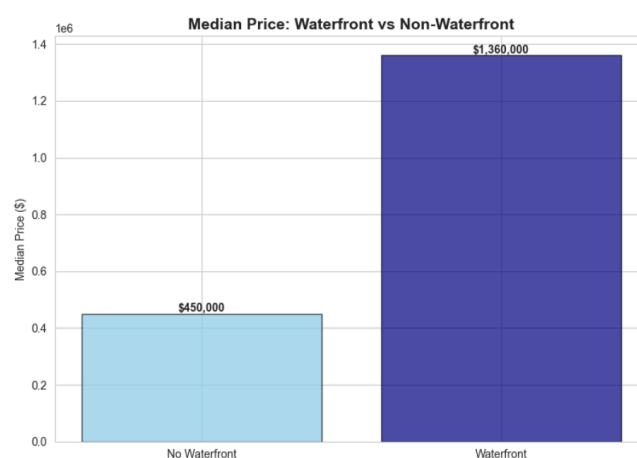


Waterfront vs Non-Waterfront Properties:

Median prices were compared for properties with and without waterfront access.

Observations:

- * Waterfront properties have a substantially higher median price.
- * Confirms that proximity to water is a major value driver.
- * The clear separation justifies treating waterfront as an important categorical feature.

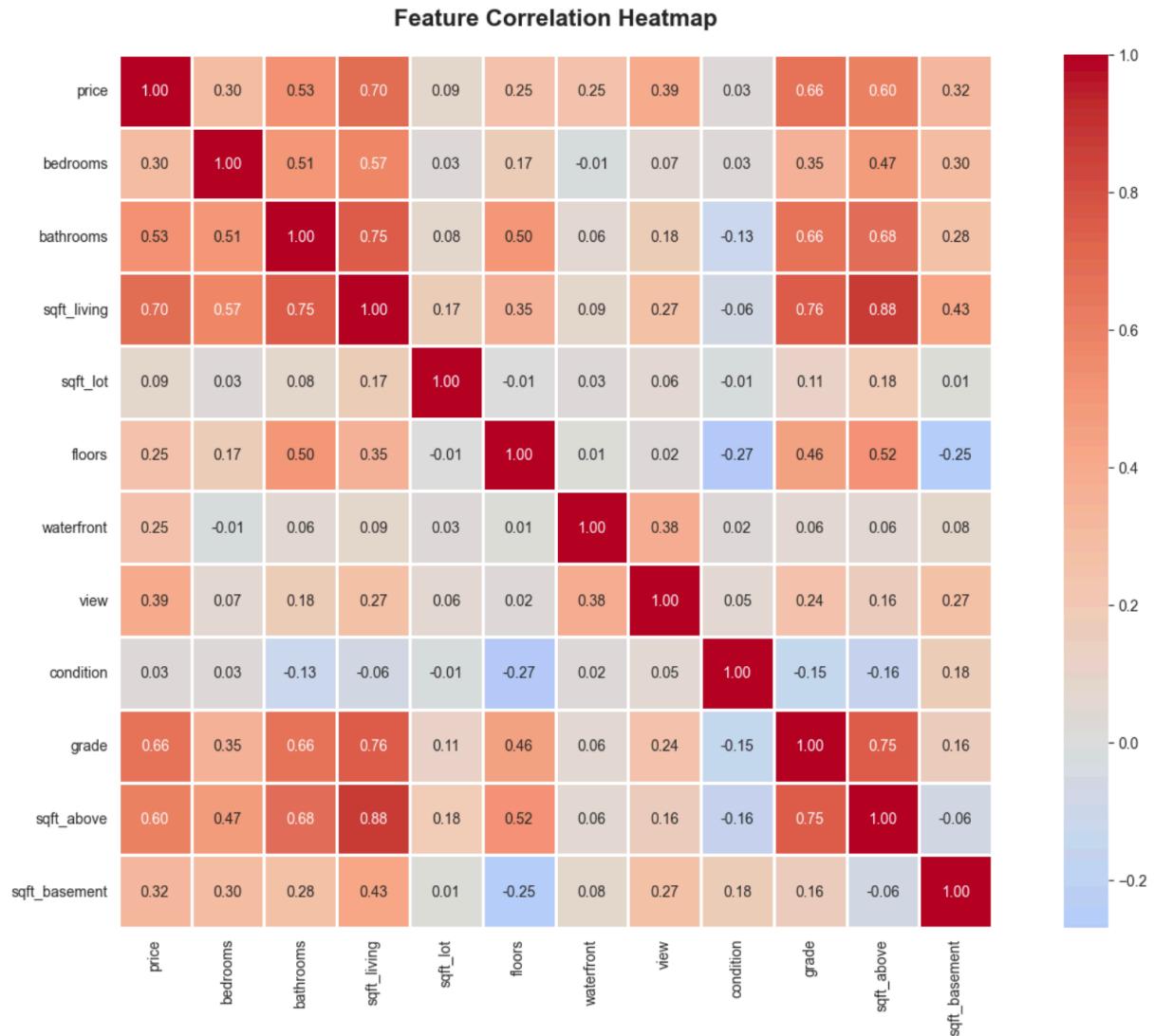


Feature Correlation Analysis:

A correlation heatmap was generated to analyse linear relationships between numerical and ordinal features.

Key observations:

- * sqft_living, grade, `bathrooms, and sqft_above show strong positive correlation with price.
- * No single feature explains price entirely — reinforcing the need for feature combination.
- * Some features are correlated with each other, indicating potential multicollinearity.



3. Feature Engineering

Feature engineering was performed to enrich the dataset with temporal, structural, relational, and geospatial information that is not explicitly available in the raw features. The goal was to improve the model's ability to capture real-world factors influencing property prices while keeping the transformations interpretable.

Temporal Feature Engineering:

The transaction date was converted into a datetime format and decomposed into multiple components.

- * Year: Captures long-term price trends and market shifts.
- * Month: Accounts for seasonal effects in the real estate market.
- * Day of Week: Encodes weekly transaction patterns.

Renovation-Based Features:

The original yr_renovated feature contains zero values indicating no renovation. To make this information more informative:

- * is_renovated: Binary indicator capturing whether a house has ever been renovated.
- * years_since_renovation:
 - * For renovated houses: difference between sale year and renovation year.
 - * For non-renovated houses: difference between sale year and construction year.

Property Age Feature:

A dedicated age feature was created to represent how old a property is at the time of sale.

house_age = (most recent sale year – year built)

Structural Consistency Validation:

To ensure data integrity, a consistency check was performed:

- * Verified whether sqft_living \approx sqft_above + sqft_basement
- * Counted inconsistent records where the difference exceeded a small tolerance.

Basement Feature Engineering:

Basements contribute differently to property valuation compared to above-ground living space.

has_basement: Binary indicator capturing basement presence.

Size-Based Ratio Features:

Raw size values were converted into normalised ratios* to improve comparability across properties.

- * living_lot_ratio: Ratio of living area to total lot size.
- * above_living_ratio: Proportion of above-ground living space.

Neighbourhood Comparison Features:

To contextualise each property relative to its neighbourhood, comparative features were created.

- * living_vs_neighbours: Ratio of property living area to average living area of nearby houses.
- * lot_vs_neighbours: Ratio of property lot size to average neighbourhood lot size.

Aggregate Room Feature:

To capture overall property capacity:

total_rooms = bedrooms + bathrooms

Geospatial Feature Engineering:

Location-based features were explicitly engineered to capture spatial pricing effects.

Distance from City Center:

- * Computed approximate distance from Seattle city centre using latitude and longitude.
- * Converted geographic distance into kilometres.

Geographic Waterfront Proximity (Heuristic):

A simplified spatial heuristic was used to estimate proximity to water bodies based on latitude and longitude ranges.

- * near_water_geo: Binary indicator for geographic closeness to waterfront areas.

4. Architecture Diagram (Tabular + Image Fusion)

Model Architecture Description:

The multimodal architecture consists of two parallel branches:

1. Tabular Data Model

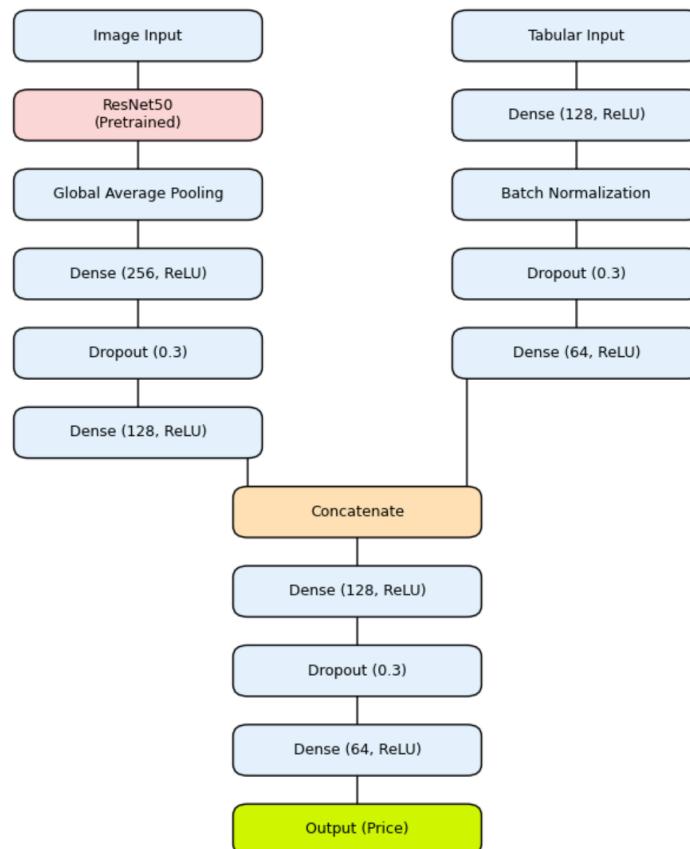
- * Input: Structured numerical features
- * Processing: Fully connected (Dense) layers
- * Output: Learned tabular feature embedding

2. Image Model

- * Input: Satellite images
- * Processing: Convolutional Neural Network (ResoNet + External layers)
- * Output: Image feature embedding

Fusion Layer

- * Concatenation of tabular and image embeddings
- * Followed by dense layers for joint learning
- * Final regression output for price prediction



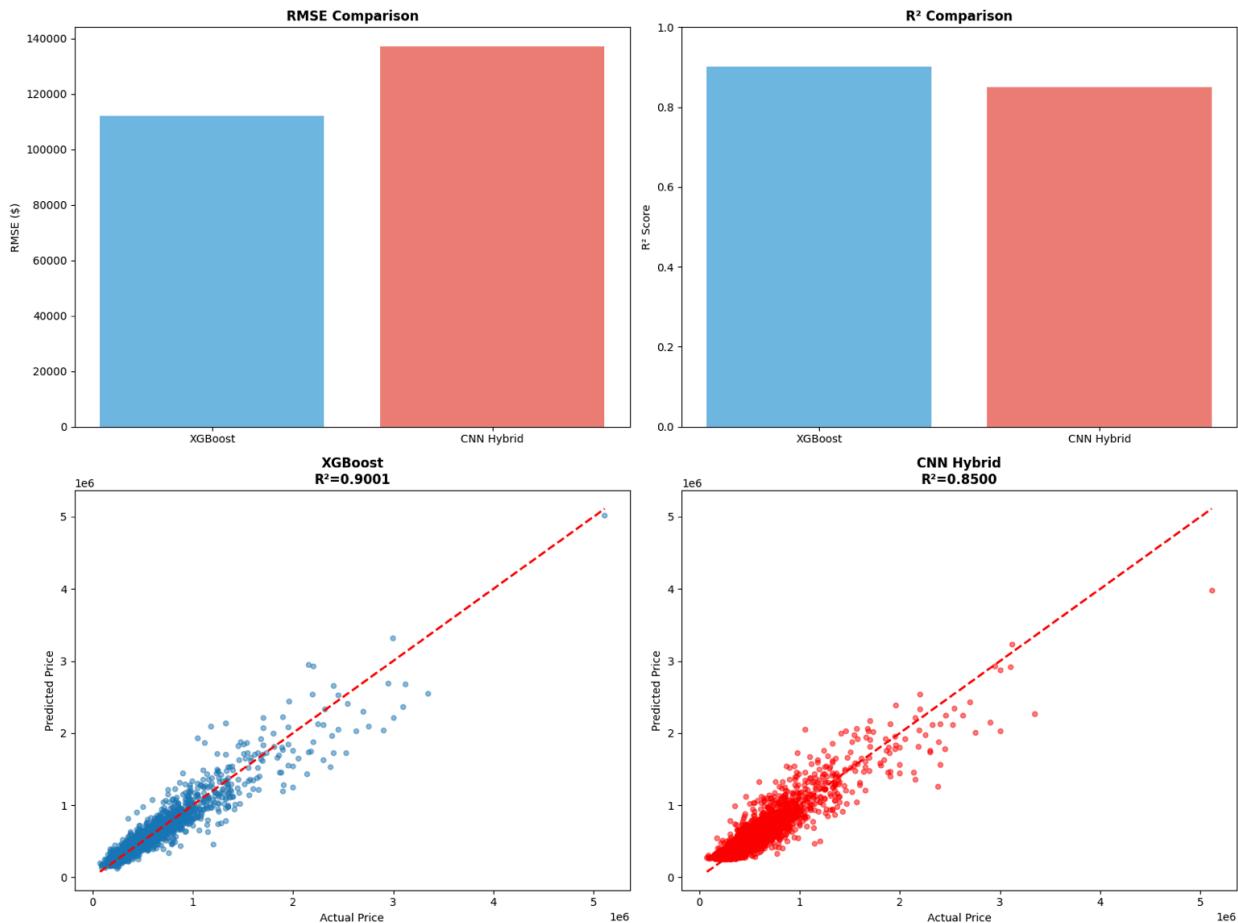
5. Results and Performance Comparison

Evaluation Setup:

Models were evaluated using consistent train–test splits and standard regression metrics.

Model Comparison:

Model	RMSE	R ²
XGBoost	111979.120339	0.900076
CNN Hybrid	137251.383498	0.850020
Best Model (RMSE):	XGBoost	
Best Model (R ²):	XGBoost	



6. Model Interpretability using Grad-CAM

While multimodal deep learning models achieve strong predictive performance, they often suffer from limited interpretability. To address this, Gradient-weighted Class Activation Mapping (Grad-CAM) was employed to visualise which regions of satellite images contribute most to the model's price predictions.

Since satellite imagery plays a critical role in the multimodal model, it is important to ensure that:

- * The CNN focuses on meaningful spatial patterns rather than noise.
- * Visual features influencing predictions align with real-world intuition.
- * The model's decisions are transparent and explainable.

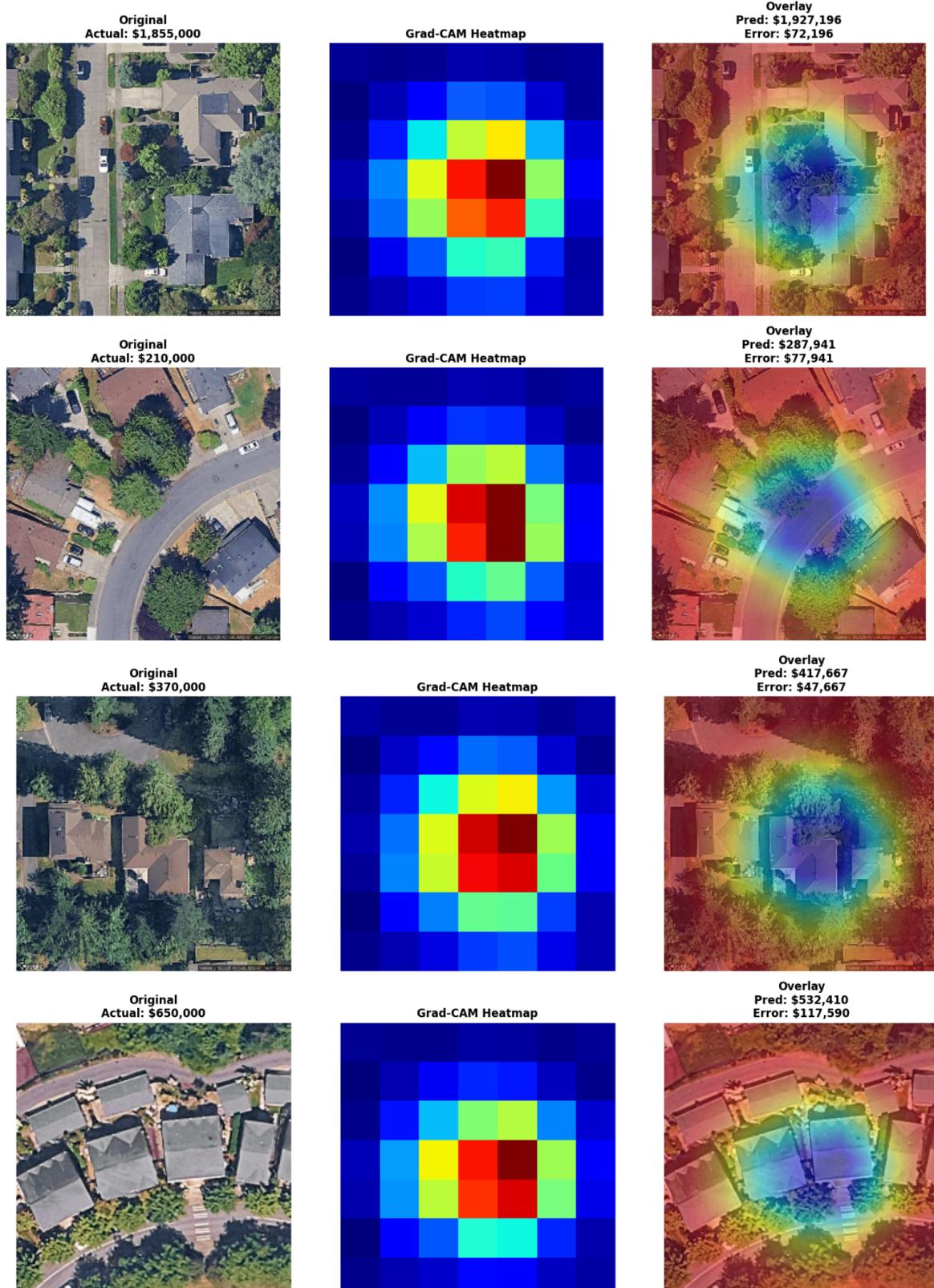
Grad-CAM provides spatial heatmaps highlighting image regions that most strongly influence the predicted output.

Grad-CAM was applied to the final convolutional layer of the CNN branch to preserve high-level semantic information.

- * The gradient of the predicted price with respect to convolutional feature maps was computed.
- * Gradients were spatially averaged to obtain importance weights.
- * Feature maps were combined using weighted summation.
- * The resulting heatmap was normalised and overlaid on the original satellite image.

For qualitative analysis:

- * Four test samples were randomly selected.
- * For each sample, the following were visualised:
 1. Original satellite image (with actual price)
 2. Grad-CAM heatmap
 3. Heatmap overlaid on the original image (with predicted price and absolute error)



Regions with High Activation

- * Open green spaces and vegetation
- * Planned residential layouts
- * Road networks and accessibility
- * Low-density housing areas

Regions with Low Activation

- * Homogeneous concrete clusters
- * Non-residential or industrial regions
- * Visually noisy or irrelevant areas

7. Conclusion

The project demonstrates that integrating satellite imagery with structured data leads to more accurate and robust real estate price prediction if trained properly with right architecture. Multimodal learning captures both quantitative attributes and qualitative neighbourhood signals, offering a scalable and practical approach for modern valuation systems.
