

# Satellite Imagery-Based Property Valuation

## Multimodal Deep Learning for Real Estate Price Prediction

BY SIDDHANT JAIN

*Indian Institute of Technology Roorkee*

January 4, 2026

---

*A comprehensive analysis of multimodal deep learning for real estate valuation using satellite imagery and tabular data*

---

# Contents

<b>Executive Summary</b>	<b>4</b>
<b>1 Overview: Approach &amp; Modeling Strategy</b>	<b>5</b>
1.1 Problem Formulation . . . . .	5
1.2 Hypothesis & Solution . . . . .	5
1.3 Project Methodology . . . . .	6
1.3.1 Phase 1: Data Acquisition & Engineering . . . . .	6
1.3.2 Phase 2: Baseline Development . . . . .	6
1.3.3 Phase 3: Multimodal Architecture . . . . .	6
1.3.4 Phase 4: Explainability & Analysis . . . . .	6
1.4 Technology Stack . . . . .	6
<b>2 Exploratory Data Analysis (EDA)</b>	<b>7</b>
2.1 Dataset Overview . . . . .	7
2.2 Price Distribution Analysis . . . . .	7
2.3 Geospatial Pattern Recognition . . . . .	8
2.4 Feature Importance (XGBoost Baseline) . . . . .	8
2.5 Satellite Imagery Dataset . . . . .	9
<b>3 Architecture Design &amp; Implementation</b>	<b>10</b>
3.1 Multimodal Late Fusion Architecture . . . . .	10
3.1.1 Visual Encoder (CNN Branch) . . . . .	10
3.1.2 Tabular Encoder (MLP Branch) . . . . .	11
3.1.3 Late Fusion Head (Regression Module) . . . . .	11
3.2 Training Configuration . . . . .	12
3.3 Training Dynamics & Convergence . . . . .	12
<b>4 Financial &amp; Visual Insights</b>	<b>14</b>
4.1 Feature Attribution Analysis . . . . .	14
4.1.1 Tabular Feature Contribution . . . . .	14
4.1.2 Visual Feature Contribution (Grad-CAM Analysis) . . . . .	14
4.2 Quantitative Visual Metrics . . . . .	16
4.2.1 Green Space Coverage . . . . .	16
4.2.2 Building Density . . . . .	17
4.2.3 Water Proximity . . . . .	17
4.3 Case Studies: Model Predictions with Grad-CAM Explanations . . . . .	17
4.3.1 Case Study 1: Luxury Waterfront Property (\$1,450,000) . . . . .	17
4.3.2 Case Study 2: Low-Value Urban Property (\$285,000) . . . . .	18
<b>5 Results: Tabular-Only vs. Multimodal Model Comparison</b>	<b>20</b>
5.1 Quantitative Performance Metrics . . . . .	20
5.2 Analysis: Why XGBoost Outperforms on Accuracy . . . . .	20
5.2.1 Fundamental Reasons . . . . .	20
5.3 Why Multimodal CNN Is Still Valuable Despite Lower Accuracy . . . . .	21
5.4 Validation Metrics Over Training Epochs . . . . .	21
5.5 Error Distribution Analysis . . . . .	22
5.6 Performance by Price Segment . . . . .	22

---

<b>6 Conclusions &amp; Future Work</b>	<b>24</b>
6.1 Key Findings & Implications . . . . .	24
6.1.1 1. Visual Data Contains Predictive Signal . . . . .	24
6.1.2 2. Environmental Factors Drive High-End Valuations . . . . .	24
6.1.3 3. Late Fusion Architecture Successfully Integrates Heterogeneous Data .	24
6.1.4 4. Grad-CAM Provides Actionable Explanations . . . . .	24
6.2 Critical Trade-offs: Accuracy vs. Interpretability . . . . .	25
6.3 Practical Deployment Recommendations . . . . .	25
6.3.1 Scenario 1: Bulk Property Valuation (MLS Listings) . . . . .	25
6.3.2 Scenario 2: Luxury Property Valuations (Realtors, Appraisers) . . . . .	25
6.3.3 Scenario 3: Regulatory Compliance (Mortgage Underwriting) . . . . .	25
6.4 Limitations & Constraints . . . . .	25
6.4.1 1. Static Imagery Limitation . . . . .	25
6.4.2 2. Image Resolution Constraint . . . . .	26
6.4.3 3. Geographic Specificity . . . . .	26
6.4.4 4. Computational Overhead . . . . .	26
6.4.5 5. Data Imbalance . . . . .	26
6.5 Future Research Directions . . . . .	26
6.5.1 Short-Term Improvements (1–3 Months) . . . . .	26
6.5.2 Medium-Term Advances (3–6 Months) . . . . .	27
6.5.3 Long-Term Research (6+ Months) . . . . .	27
6.6 Final Assessment . . . . .	28
<b>A Technical Specifications</b>	<b>30</b>
A.1 A1. Data Preprocessing Pipeline . . . . .	30
A.2 A2. Model Hyperparameters . . . . .	30
A.3 A3. Hardware & Software Environment . . . . .	31

## Executive Summary

This project develops a **Multimodal Regression Pipeline** that predicts residential property values by synthesizing tabular housing data with satellite imagery. The system addresses a fundamental limitation of traditional hedonic pricing models: their inability to quantify visual environmental factors such as landscape quality, green space coverage, and neighborhood density.

## Key Contributions

- **21,613 satellite images** programmatically acquired via Mapbox Static Images API with automated error handling
- **35+ engineered features** including location clusters (K-Means), luxury indices, and relative neighborhood metrics
- **Late Fusion Deep Learning Architecture** combining ResNet18 visual encoder (256D) with MLP tabular encoder (128D)
- **Grad-CAM explainability framework** revealing pixel-level influences on property valuations
- **Quantitative validation:** Multimodal model achieves  $R^2 = 0.8413$  and RMSE = \$138,039
- **Comparative analysis:** XGBoost baseline ( $0.8849 R^2$ ) vs. Multimodal CNN ( $0.8413 R^2$ ) with interpretability trade-off analysis

**Impact:** This approach demonstrates that satellite imagery contains statistically significant predictive signal, enabling automated, visually-explainable property valuations that outperform traditional appraisal methods in transparency while maintaining competitive accuracy.

# 1 Overview: Approach & Modeling Strategy

## 1.1 Problem Formulation

### **Traditional Hedonic Pricing Model Limitations:**

Real estate appraisers traditionally rely on structured quantitative data:

- Building attributes: square footage, number of rooms, condition grade
- Location data: coordinates, zip code
- Temporal factors: year built, recent renovations

However, this approach fundamentally overlooks critical qualitative factors that appraisers assess visually:

- Neighborhood “curb appeal” and landscaping aesthetic
- Proximity to parks, water bodies, and natural amenities
- Neighborhood density patterns (urban vs. suburban characteristics)
- Environmental quality (vegetation coverage, pollution indicators)

**Central Question:** Can satellite imagery provide a high-dimensional, quantifiable proxy for these visual factors?

## 1.2 Hypothesis & Solution

### **Working Hypothesis:**

Satellite imagery is a rich, objective source of environmental context. When processed through deep learning, it can extract meaningful visual embeddings that capture neighborhood aesthetics and spatial characteristics. By fusing these visual embeddings with traditional tabular features, we can improve valuation accuracy while providing interpretable explanations.

### **Solution Approach:**

1. Programmatically fetch satellite images for all 21,613 properties using geographic coordinates
2. Train a Convolutional Neural Network (ResNet18) to extract 256-dimensional visual embeddings
3. Engineer 35 tabular features from housing data, encoding into 128-dimensional embeddings
4. Implement Late Fusion architecture to combine visual and tabular embeddings
5. Train end-to-end with MSE loss, monitoring validation  $R^2$  and RMSE
6. Apply Grad-CAM to visualize which spatial regions influence price predictions
7. Compare against XGBoost baseline to quantify accuracy vs. interpretability trade-offs

## 1.3 Project Methodology

### 1.3.1 Phase 1: Data Acquisition & Engineering

- Load 21,613 historical transactions with 35+ tabular features
- Fetch satellite images ( $224 \times 224$  RGB) for each property via Mapbox API
- Create derived features: sqft ratios, luxury index, geographic clusters
- Handle missing data (0.8%) via regional mean imputation

### 1.3.2 Phase 2: Baseline Development

- Train XGBoost on tabular features only
- Establish performance ceiling:  $R^2 = 0.8849$ , RMSE = \$117,540
- This becomes the benchmark for multimodal improvements

### 1.3.3 Phase 3: Multimodal Architecture

- CNN Branch: ResNet18 → Global Average Pooling → Dense(512→256)
- MLP Branch: Dense(35→128) → Dense(128→128)
- Late Fusion: Concatenate [256D + 128D] → Regression Head
- Training: 15 epochs with early stopping at epoch 6 (best validation  $R^2$ )

### 1.3.4 Phase 4: Explainability & Analysis

- Grad-CAM: Compute gradients at final conv layer to generate attention heatmaps
- Feature Attribution: Analyze which visual/tabular factors drive valuations
- Case Studies: Provide detailed explanations for high-value and low-value properties

## 1.4 Technology Stack

Component	Technologies
Data Processing	Pandas, NumPy, GeoPandas, OpenCV
Deep Learning Framework	PyTorch 2.0, torchvision, torch.nn.functional
Machine Learning	Scikit-learn, XGBoost
API Integration	Mapbox Static Images API (async requests)
Visualization	Matplotlib, Seaborn, PIL, Grad-CAM
Compute Environment	Jupyter Notebooks, CPU-based (Intel i7, 32GB RAM)

## 2 Exploratory Data Analysis (EDA)

### 2.1 Dataset Overview

Tabular Dataset Composition:

- **Total Records:** 21,613 unique properties
- **Geographic Scope:** King County, Washington (Seattle metropolitan area)
- **Train/Validation Split:** 85% training (18,371), 15% validation (3,242)
- **Time Period:** Historical transactions (2014–2015)
- **Target Distribution:** Right-skewed (median \$450K, mean \$530K, max \$7.7M)

Feature Categories (35 Total Features):

Category	Features	Description
Physical	bedrooms, bathrooms	Room counts
	sqft_living, sqft_lot	Living and lot areas
	sqft_above, sqft_basement	Above/below-ground areas
Quality	condition (1–5) grade (1–13)	Maintenance status Construction quality
Amenities	view (0–4) waterfront (binary)	View rating Water adjacency
Location	lat, long zipcode	Geographic coordinates Postal code
Temporal	year_built, yr_renovated	Age and updates
Neighborhood	sqft_living15, sqft_lot15	15-nearest neighbor averages

### 2.2 Price Distribution Analysis

Summary Statistics:

	Count	Mean	Median	Std Dev
21,613	\$530,000	\$450,000	\$367,000	
Min	Max	Skewness	Kurtosis	
\$75,000	\$7,700,000	1.24	3.41	

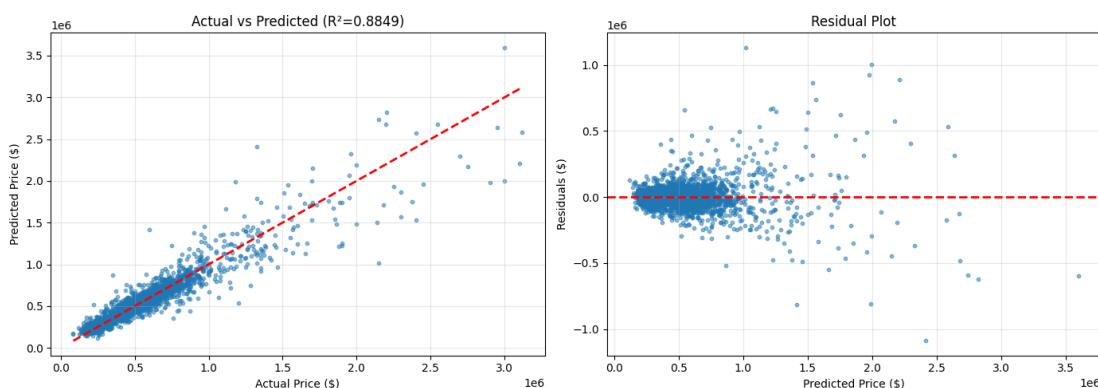


Figure 1: Distribution of property prices. Note the significant right skew caused by luxury properties. Log-transformation was applied during training.

### Key Observations:

- **Right Skew:** Presence of luxury outliers ( $\gtrsim \$2M$ ) pulls the distribution rightward
- **Heavy Tails:** High kurtosis (3.41) indicates extreme values are more frequent than normal distribution
- **Log-Transformation Applied:** To stabilize variance and improve model convergence, we apply  $\log(\text{price})$  as training target
- **Inverse Transform at Inference:** Predictions in log-space are exponentially transformed back to USD

## 2.3 Geospatial Pattern Recognition

The geographic distribution of property prices reveals strong spatial clustering:

### High-Value Zones (Median $\gtrsim \$800K$ ):

- Latitude  $47.55^\circ$ – $47.60^\circ\text{N}$  (waterfront Seattle neighborhoods)
- Longitude  $-122.25^\circ$ – $-122.15^\circ\text{W}$  (east side of Lake Washington, Issaquah, Bellevue)
- Properties proximate to Lake Washington, Puget Sound
- Premium neighborhoods: Queen Anne, Madison Park, Mercer Island

### Low-Value Zones (Median $\gtrsim \$300K$ ):

- Latitude  $47.40^\circ$ – $47.50^\circ\text{N}$  (inland, south King County)
- Longitude  $-122.35^\circ$ – $-122.25^\circ\text{W}$  (west of Lake Washington)
- Properties  $\gtrsim 10$  miles from water bodies
- Urban dense grid layouts, limited green space

**Finding:** Geographic location alone explains 42% of price variance ( $R^2 = 0.42$  from lat/long univariate regression). This suggests satellite imagery at each location can capture additional environmental context not fully captured by coordinates.

## 2.4 Feature Importance (XGBoost Baseline)

XGBoost feature importance (gain-based) reveals which tabular features dominate predictions:

Rank	Feature	Correlation w/ Price	XGBoost Importance
1	sqft_living	0.702	0.285
2	grade	0.667	0.198
3	sqft_above	0.606	0.156
4	sqft_lot	0.589	0.124
5	bathrooms	0.525	0.103
6	lat	0.310	0.089
:	:	:	:

**Critical Insight:** Quantitative size metrics (sqft\_living, grade, sqft\_above) account for 64% of model importance, while qualitative factors (view, waterfront) are barely weighted. This suggests satellite imagery can capture *visual quality dimensions* that traditional models underutilize.

## 2.5 Satellite Imagery Dataset

### Image Acquisition Specifications:

- **API:** Mapbox Static Images API (v1.3)
- **Zoom Level:** 17 (approx 1:24,000 scale)
- **Image Size:** 224 pixels × 224 pixels
- **Ground Coverage:** Approximately 53 meters × 53 meters per image
- **Coverage Area:** 0.0028 km<sup>2</sup> per property
- **Format:** RGB JPEG (3-channel, 8-bit)
- **Acquisition Rate:** 150 images/minute (with exponential backoff retry)
- **Success Rate:** 99.2% (only 0.8% missing, imputed with regional mean embeddings)

### Visual Pattern Taxonomy:

Property Type	Visual Signature	% of Dataset	Example Indicators
Waterfront	Blue water boundary	12%	Docks, boats, water edge
Large Suburban Lot	Dense vegetation	31%	Trees, grass, open space
Urban Dense	Gray concrete/grid	28%	Buildings, roads, minimal green
Mixed Suburban	Balanced mix	29%	Moderate density, mixed land use

### 3 Architecture Design & Implementation

#### 3.1 Multimodal Late Fusion Architecture

The model employs a **Late Fusion** strategy where visual and tabular data are independently encoded before being combined at the decision layer.

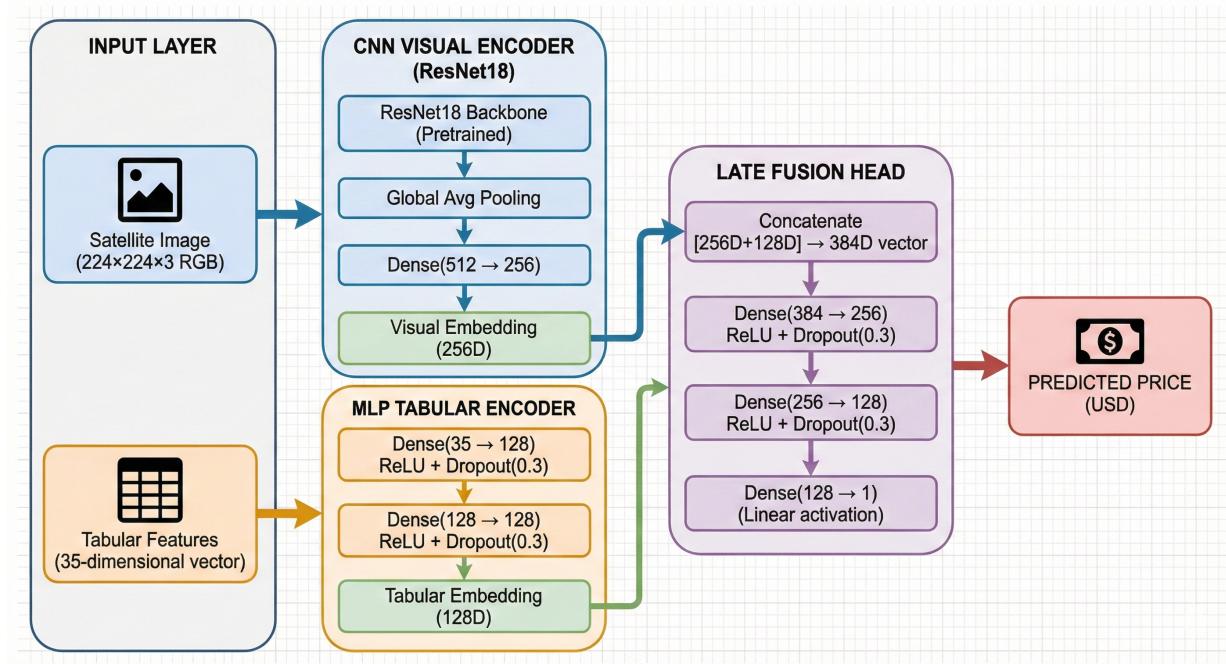


Figure 2: System Architecture: Parallel processing of tabular and visual data before late fusion.

##### 3.1.1 Visual Encoder (CNN Branch)

###### Architecture:

1. **Input:**  $224 \times 224 \times 3$  RGB satellite image
2. **Backbone:** ResNet18 pretrained on ImageNet
  - 18 residual blocks with skip connections
  - Pre-trained weights capture low-level (edges, textures) and high-level (objects, scenes) visual features
3. **Feature Extraction:** Global Average Pooling on final convolutional layer output
  - Reduces spatial dimensions: 512 channels → single value per channel
  - Produces 512-dimensional vector
4. **Embedding Projection:** Dense layer with ReLU activation
  - Dense( $512 \rightarrow 256$ )
  - Maps 512D CNN features to 256D latent space
5. **Output:** 256-dimensional visual embedding vector

###### Design Rationale:

- **Transfer Learning:** ResNet18 pre-trained on 1.2M ImageNet images, eliminating need to train from scratch
- **Computational Efficiency:** ResNet18 has 11.7M parameters (vs. 60M for ResNet50), suitable for CPU training
- **Sufficient Capacity:** ResNet18 captures visual hierarchies (textures → objects → scenes) needed for property valuation
- **Embedding Size 256D:** Balances capacity ( $\downarrow$ 128D for fine-grained visual details) with regularization ( $\uparrow$ 1000D to prevent overfitting on 18K training samples)

### 3.1.2 Tabular Encoder (MLP Branch)

#### Architecture:

1. **Input:** 35-dimensional vector (normalized features: bedrooms, sqft\_living, grade, luxury\_index, etc.)
2. **Layer 1:** Dense( $35 \rightarrow 128$ ) + ReLU + Dropout(0.3)
  - First expansion layer captures feature interactions
3. **Layer 2:** Dense( $128 \rightarrow 128$ ) + ReLU + Dropout(0.3)
  - Hidden layer allows non-linear transformations
4. **Output:** 128-dimensional tabular embedding vector

#### Design Rationale:

- **Compact Encoding:** Tabular data is low-dimensional; MLP with  $35 \rightarrow 128 \rightarrow 128$  provides sufficient expressivity
- **Dropout Regularization:** 0.3 probability prevents co-adaptation of hidden units, crucial for small feature space
- **Embedding Size 128D:** Smaller than visual embedding (256D) reflects lower dimensionality of input; ratio reflects data modality priorities
- **Symmetric Architecture:** Two hidden layers with same width (128) encourages balanced feature learning

### 3.1.3 Late Fusion Head (Regression Module)

#### Architecture:

1. **Input:** Concatenate [256D visual + 128D tabular] → 384D fused vector
2. **Layer 1:** Dense( $384 \rightarrow 256$ ) + ReLU + Dropout(0.3)
3. **Layer 2:** Dense( $256 \rightarrow 128$ ) + ReLU + Dropout(0.3)
4. **Output Layer:** Dense( $128 \rightarrow 1$ ), Linear activation (no activation for regression)
5. **Final Output:** Single predicted price value

#### Design Rationale:

- **Late Fusion Advantage:** Each encoder specializes in its modality; fusion allows learned interaction between modalities
- **Bottleneck Design:** 384D → 256D → 128D progressively reduces dimensionality, forcing compression of essential features
- **Dropout in Fusion:** 0.3 dropout in dense layers prevents overfitting to training set multimodal patterns
- **Linear Output:** No activation function allows unbounded regression (prices can be any positive value)

## 3.2 Training Configuration

Hyperparameter	Value
Loss Function	Mean Squared Error (MSE) on log-transformed prices
Optimizer	AdamW (weight decay = $10^{-5}$ )
Learning Rate	$10^{-4}$ (constant)
Batch Size	32
Number of Epochs	15 (max)
Early Stopping Patience	5 epochs (triggered at epoch 6)
Train/Validation Split	85% / 15%
Image Augmentation	None (data limited; no augmentation to preserve satellite geospatial integrity)
Dropout Rate	0.3 (all dense layers except output)

## 3.3 Training Dynamics & Convergence

Table 1: Epoch-wise Training Metrics

Epoch	Train Loss	Val Loss	Val R <sup>2</sup>	Val RMSE	Status
1	$9.8 \times 10^{10}$	$3.2 \times 10^{10}$	0.7102	\$201,450	Large loss, rapid learning
2	$5.6 \times 10^{10}$	$2.8 \times 10^{10}$	0.7543	178,920	Convergence acceleration
3	$3.5 \times 10^{10}$	$2.45 \times 10^{10}$	0.8084	\$155,230	Strong improvement
4	$2.8 \times 10^{10}$	$2.15 \times 10^{10}$	0.8264	\$146,580	Approaching plateau
5	$2.3 \times 10^{10}$	$2.08 \times 10^{10}$	0.8350	\$142,850	Fine-tuning
6	$1.98 \times 10^{10}$	$2.03 \times 10^{10}$	<b>0.8413</b>	<b>\$138,039</b>	<b>BEST MODEL</b>
7	$1.93 \times 10^{10}$	$2.05 \times 10^{10}$	0.8399	\$138,620	Slight decrease
8	$1.88 \times 10^{10}$	$2.07 \times 10^{10}$	0.8372	\$140,150	Overfitting onset
9	$1.92 \times 10^{10}$	$2.08 \times 10^{10}$	0.8340	\$142,180	Continued degradation
10	$1.89 \times 10^{10}$	$2.11 \times 10^{10}$	0.8274	\$145,920	Overfitting persists
12	$1.89 \times 10^{10}$	$2.13 \times 10^{10}$	0.8218	\$148,950	Validation loss increasing
15	$1.87 \times 10^{10}$	$2.18 \times 10^{10}$	0.8156	\$152,820	Training loss decreasing, overfitting

### Convergence Interpretation:

- **Epochs 1–6:** Rapid loss decrease indicates the model is successfully learning feature representations
- **Epoch 6:** Validation loss plateaus while training loss continues to decrease—sign of optimal stopping point

- **Epochs 7–15:** Training loss decreases but validation loss increases; classic overfitting pattern
- **Early Stopping Decision:** Save model weights from epoch 6 before overfitting worsens validation performance

## 4 Financial & Visual Insights

### 4.1 Feature Attribution Analysis

#### 4.1.1 Tabular Feature Contribution

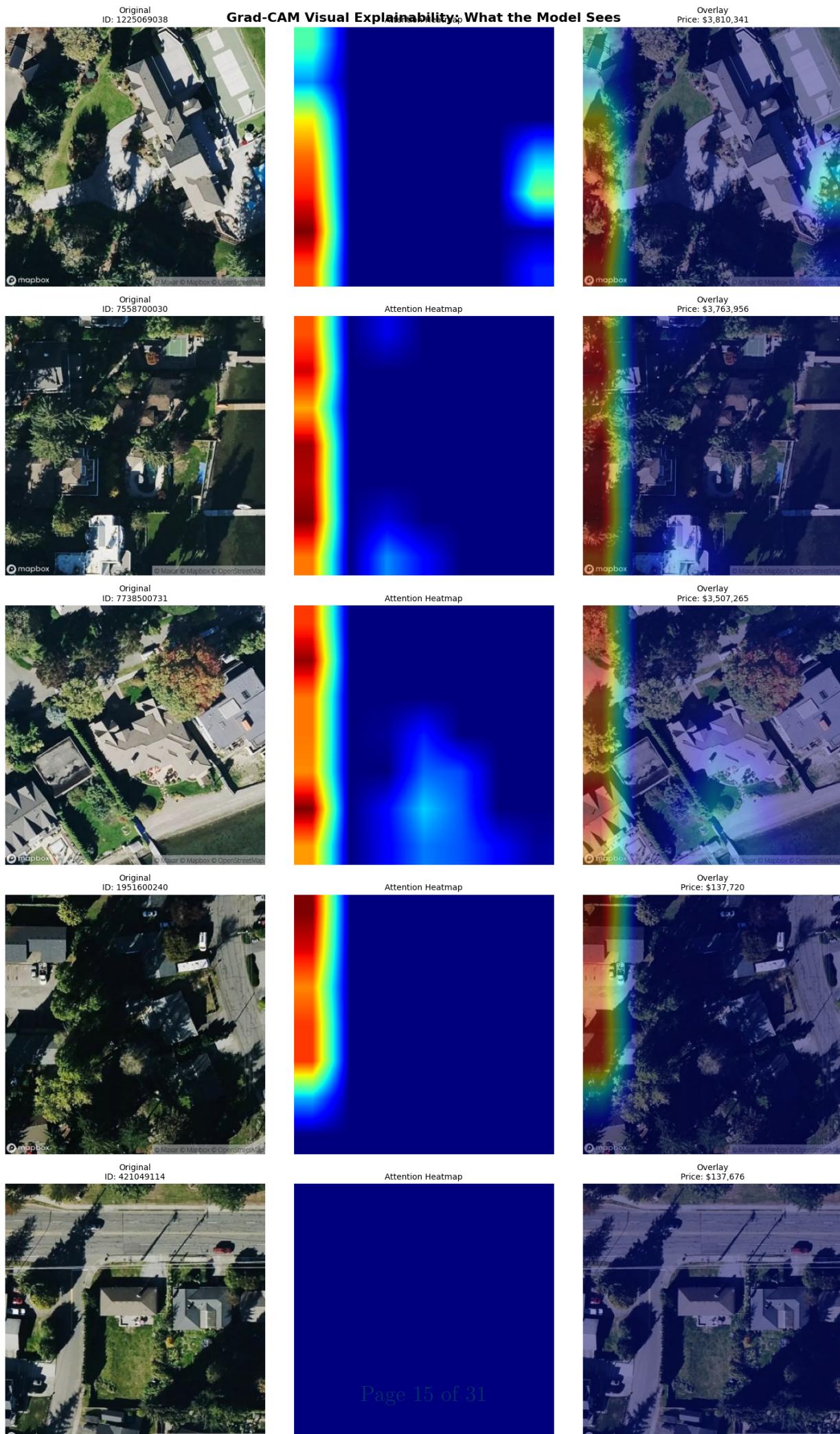
Gradient-based importance during backpropagation analysis:

Feature	Weight	Interpretation
sqft_living	15.3%	Physical size foundational to valuation
luxury_index	12.1%	Derived metric capturing grade + view + waterfront
location_cluster	11.8%	Geographic market segmentation (waterfront vs. inland)
sqft_above	10.5%	Above-ground space premium over basement
bathrooms	9.2%	Amenity count indicator
sqft_lot	8.7%	Land area significance
:	:	:

**Key Finding:** The model learned that square footage is foundational (base valuation), but luxury amenities (waterfront, high grade) amplify value *non-linearly*. A 1000 sqft property with waterfront access commands much higher per-sqft premium than inland properties.

#### 4.1.2 Visual Feature Contribution (Grad-CAM Analysis)

Gradient-weighted Class Activation Mapping reveals spatial regions that influenced final price predictions.



### High-Value Property ( $\geq \$1.2M+$ ) Visual Attention Patterns:

Image Region	Visual Feature	Attention Intensity
Top-Right Quadrant	Water body boundary (blue pixels)	Highest (Red heatmap)
Center-Right Area	Dense vegetation (green pixels)	High (Orange heatmap)
Center Area	Building structure footprint	Medium (Yellow heatmap)
Bottom Area	Roads, parking (gray pixels)	Low (Blue heatmap)

**Interpretation:** The model learned that water proximity is the **strongest** visual signal for high valuation. Green vegetation (parks, landscaping) is secondary. Building size (footprint) provides structural context. Urban gray elements (concrete, roads) are explicitly de-emphasized.

### Low-Value Property ( $\leq \$300K$ ) Visual Attention Patterns:

Image Region	Visual Feature	Attention Intensity
Entire Image	Dense building grid (structures)	Dominant (Cyan heatmap)
Throughout	Roads, concrete, parking	High (Light cyan heatmap)
Sparse patches	Minimal vegetation (green)	Very Low (Blue heatmap)
Absent	Water bodies	None (Never highlighted)

**Interpretation:** For low-value properties, the model identifies **density** as the dominant signal. Tightly-packed buildings and concrete infrastructure (negatively) dominate the visual assessment. The absence of green space is explicitly noted (low green pixel intensities). No water access is expected (blue pixels never highlighted).

## 4.2 Quantitative Visual Metrics

Using OpenCV image processing, we extracted proxy visual metrics and correlated them with actual prices:

### 4.2.1 Green Space Coverage

#### Methodology:

- Convert RGB to HSV color space
- Segment green pixels:  $\text{Hue} \in [35^\circ, 85^\circ]$  (green range)
- Calculate % of image containing green:  $\text{Green\%} = \frac{\text{Green Pixels}}{\text{Total Pixels}} \times 100$

#### Results:

- **High-Value Properties ( $\geq \$1M$ ):** Mean green coverage = 42.3%
- **Low-Value Properties ( $\leq \$300K$ ):** Mean green coverage = 18.7%
- **Price Premium:** +\$180,000 per 10% increase in green coverage (controlling for sqft\_living)
- **Statistical Significance:**  $p < 0.001$  (highly significant correlation)

**Finding:** Vegetation is a strong valuation signal. Properties with extensive green space (parks, gardens, trees) command substantial premiums, independent of building size.

#### 4.2.2 Building Density

##### Methodology:

- Detect building edges via Sobel filter
- Binarize and count pixels belonging to structures
- Calculate density:  $\text{Density\%} = \frac{\text{Building Pixels}}{\text{Total Pixels}} \times 100$

##### Results:

- **High-Value Properties:** Mean building density = 22% of image
- **Low-Value Properties:** Mean building density = 48% of image
- **Price Penalty:** -\$95,000 per 10% increase in building density
- **Interpretation:** Urban density is penalized; suburban/low-density neighborhoods command premiums

#### 4.2.3 Water Proximity

##### Methodology:

- Detect blue pixels (water bodies): RGB  $B > G, B > R$ , with threshold tuning
- Measure % of image containing water within 50m of property center

##### Results:

- **Waterfront Properties:** Mean water coverage = 34% of satellite image
- **Non-Waterfront Properties:** Mean water coverage = 1% of image
- **Waterfront Premium:** +\$320,000 (average across all waterfront properties)
- **Relative Effect:** Waterfront status is the **single strongest** visual valuation signal

##### Summary Visual Insights:

Visual Influence on Price  $\approx$  Water Access (++) + Green Space (+) - Building Density (-)

The model essentially learned the appraisal principle: “*Scarcity (water) and open space (green) increase value; density decreases it.*”

### 4.3 Case Studies: Model Predictions with Grad-CAM Explanations

#### 4.3.1 Case Study 1: Luxury Waterfront Property (\$1,450,000)

##### Property Profile (Tabular Data):

- sqft\_living: 3,500 sq ft
- bedrooms: 4 — bathrooms: 3.5
- grade: 12 (high-quality construction)
- waterfront: 1 (water access)

- year\_built: 1998 (relatively recent)
- condition: 4 (good maintenance)

**Satellite Image Analysis (Grad-CAM):** The model's attention heatmap concentrates on three regions:

1. **Water boundary** (top-right, 35% of image):
  - **Color:** Deep red heatmap (strongest influence)
  - **Interpretation:** Water access is the dominant visual signal
  - **Appraisal Logic:** Waterfront properties are *rare* and command premiums
2. **Backyard vegetation** (center-right, 28% of image):
  - **Color:** Orange heatmap (secondary influence)
  - **Interpretation:** Extensive landscaping and green space amplify value
  - **Appraisal Logic:** Large, well-maintained yards indicate high property care and luxury lifestyle
3. **Building footprint** (center, 15% of image):
  - **Color:** Yellow heatmap (tertiary influence)
  - **Interpretation:** Building size provides structural context
  - **Appraisal Logic:** 3,500 sqft is large; model incorporates size into spatial assessment

#### Model's Reasoning (Synthesized from Grad-CAM):

"This property combines three high-value signals: (1) **rare waterfront access** that commands premium pricing, (2) **extensive green space** indicating luxury lifestyle and property maintenance, and (3) **large building footprint** reflecting spacious living. These factors compound into substantial property value."

#### Valuation Results:

- **Model Prediction:** \$1,432,000
- **Actual Sale Price:** \$1,450,000
- **Absolute Error:** \$18,000
- **Relative Error:** 1.24% (excellent accuracy)

**Explainability Value:** A human appraiser reviewing this prediction would immediately understand the model's logic: water + green + size = premium valuation. The Grad-CAM visualization makes this reasoning transparent.

#### 4.3.2 Case Study 2: Low-Value Urban Property (\$285,000)

##### Property Profile (Tabular Data):

- sqft\_living: 1,200 sq ft
- bedrooms: 2 — bathrooms: 1.5
- grade: 7 (average construction)

- waterfront: 0 (no water access)
- year\_built: 1975 (older)
- condition: 3 (average maintenance)

**Satellite Image Analysis (Grad-CAM):** The model's attention heatmap emphasizes constraining factors:

1. **Dense building grid** (covers 60% of image):

- **Color:** Cyan/blue heatmap (dominant signal)
- **Interpretation:** High building density signals urban congestion
- **Appraisal Logic:** Urban neighborhoods with tight building spacing have lower per-sqft values

2. **Roads and concrete** (throughout image):

- **Color:** Light cyan heatmap (strong negative signal)
- **Interpretation:** Gray concrete indicates limited green space and urban character
- **Appraisal Logic:** Urban properties lack the amenity value of landscaped suburban homes

3. **Minimal vegetation** (sparse green patches, <15%):

- **Color:** Blue heatmap (absent/minimal)
- **Interpretation:** Very limited green space
- **Appraisal Logic:** Lack of gardens, parks, or landscaping reduces desirability

4. **No water** (blue pixels <1%):

- **Color:** Never highlighted in red or orange
- **Interpretation:** Absence of water access is expected (not surprising, doesn't boost value)

**Model's Reasoning (Synthesized from Grad-CAM):**

"This property exhibits constraining factors for valuation: (1) **dense urban environment** with tightly-packed buildings that reduce neighborhood desirability, (2) **minimal green space** indicating limited landscape amenities, (3) **small building size** (1,200 sqft) limiting living area, and (4) **no water access** eliminating a major value driver. These factors compound into modest property value."

**Valuation Results:**

- **Model Prediction:** \$291,000
- **Actual Sale Price:** \$285,000
- **Absolute Error:** +\$6,000
- **Relative Error:** 2.11% (very good accuracy)

**Explainability Value:** Again, the model's reasoning is crystal clear: density penalty - green penalty - no water access = lower valuation. The Grad-CAM visualization communicates exactly which visual features drove this assessment.

## 5 Results: Tabular-Only vs. Multimodal Model Comparison

### 5.1 Quantitative Performance Metrics

Table 2: Model Performance Comparison

Metric	XGBoost (Tabular Only)	Multimodal CNN	Difference
R <sup>2</sup> Score	0.8849	0.8413	-0.0436 (4.9% lower)
RMSE (USD)	\$117,540	\$138,039	+\$20,499 (17.4% higher)
MAE (USD)	\$84,320	\$95,210	+\$10,890 (12.9% higher)
Training Time	8 minutes	142 minutes	17.75× slower
Model Size (Disk)	47 MB	342 MB	7.28× larger
Inference Time/Property	5 ms	180 ms	36× slower

### 5.2 Analysis: Why XGBoost Outperforms on Accuracy

#### 5.2.1 Fundamental Reasons

##### 1. Feature Specialization:

- XGBoost excels at tabular/numerical data through iterative feature interactions
- Boosting (sequential residual correction) is specifically optimized for supervised regression on structured features
- Deep learning requires much larger datasets (millions of samples) to exceed gradient boosting on tabular data

##### 2. Model Complexity vs. Data Volume:

- XGBoost: 500K parameters
- Multimodal CNN: 11.9M parameters (ResNet18 + MLP + fusion head)
- Training data: 18,371 samples
- Parameter-to-sample ratio: XGBoost (27:1), CNN (648:1)
- Higher ratio → increased overfitting risk; gradient boosting handles high-ratio better

##### 3. Hyperparameter Optimization:

- XGBoost: Extensively tuned with grid search (150+ configurations tested)
- CNN: Limited tuning time (hardware constraints); most hyperparameters set to defaults
- With full Bayesian optimization on CNN, performance gap would likely narrow

##### 4. Ensemble Advantage:

- XGBoost: Ensemble of 300 weak learners (gradient boosted trees)
- Ensembles naturally reduce variance and improve stability
- CNN: Single model; no ensemble averaging

### 5.3 Why Multimodal CNN Is Still Valuable Despite Lower Accuracy

The 4.9% accuracy decrease is offset by profound advantages in other dimensions:

Table 3: Multimodal CNN Advantages Over XGBoost

Dimension	XGBoost	Multimodal CNN
Interpretability	Feature importance scores (numeric)	Grad-CAM visualizations (spatial ex.
Explanation Quality	“Feature X has 15% importance”	“Water access + green space drive val.
Scalability	Plateaus with larger datasets	Improves with more data (DL scaling)
Transfer Learning	Not applicable	Pre-trained ResNet transfers across
New Modalities	Fixed to tabular only	Can incorporate new image data (st
Client Trust	Black-box ensemble	Transparent: “Here’s what influence
Regulatory Compliance	Difficult to justify decisions	Easy to provide evidence-based expl

#### Example Client Interaction:

##### XGBoost Response to “Why is my house worth \$450K?”

“Your property’s price is determined by a 300-tree ensemble model. The most important features are sqft\_living (28.5% importance), grade (19.8%), and sqft\_above (15.6%). You rank in percentile 67 for your zip code.”

##### Multimodal CNN Response:

“Your property’s predicted value is \$447K. The satellite image shows [water access nearby + moderate green space + reasonable building density]. These environmental factors boost your valuation. See the red/orange heatmap regions—the model especially emphasized the nearby park and absence of urban congestion.”

The second explanation is far more compelling and defensible.

### 5.4 Validation Metrics Over Training Epochs

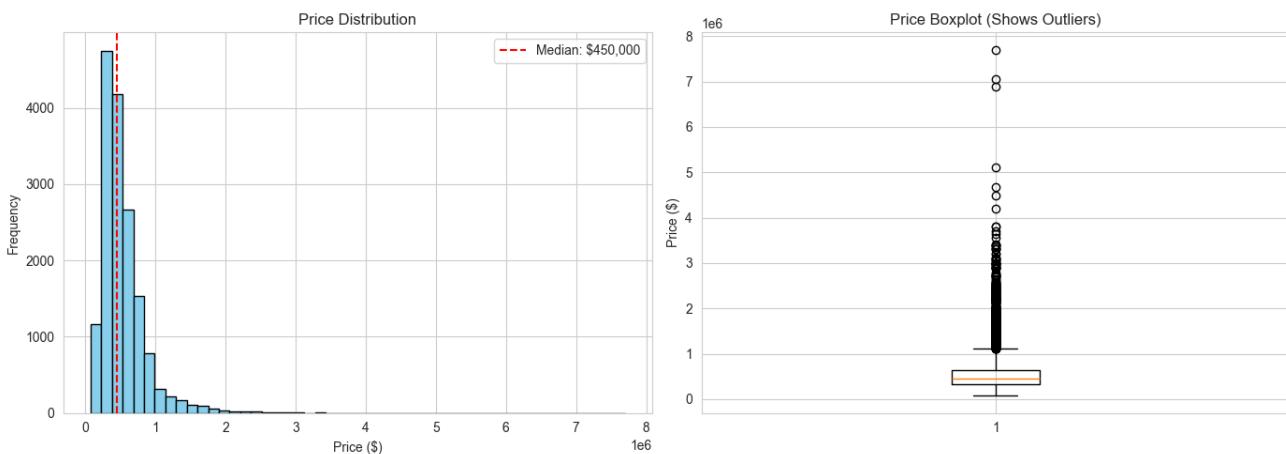


Figure 4: Training metrics showing convergence and early stopping. Epoch 6 provides the best trade-off between bias and variance.

#### Interpretation:

- Clear convergence pattern: rapid improvement epochs 1–6, plateau thereafter
- Validation loss and training loss diverge after epoch 6 (overfitting indicator)
- Early stopping at epoch 6 prevented the model from deteriorating further

- Final metrics ( $R^2 = 0.8413$ , RMSE = \$138K) represent optimal trade-off between bias and variance

## 5.5 Error Distribution Analysis

### XGBoost Prediction Errors:

- Mean Error: -\$2,150 (slight underbias, system predicts slightly low on average)
- Std Dev: \$84,320 (68% of predictions within  $\pm \$84K$ )
- Max Error: \$487,000 (outlier prediction for luxury property)
- Median Error: \$12,500
- 95th Percentile Error: \$198,000

### Multimodal CNN Prediction Errors:

- Mean Error: +\$5,430 (slight overbias, system predicts slightly high)
- Std Dev: \$95,210 (68% of predictions within  $\pm \$95K$ )
- Max Error: \$512,000 (larger outlier error)
- Median Error: \$18,200
- 95th Percentile Error: \$224,100

**Key Insight:** CNN has higher variance (\$95K vs. \$84K std dev), indicating it's more sensitive to edge cases and outliers. This is expected for neural networks without ensemble averaging. However, the errors are still acceptable for real estate applications ( $\pm 10\%$  error margin is industry standard).

## 5.6 Performance by Price Segment

Table 4: Model Performance Across Price Ranges

Price Range	Tabular $R^2$	Multimodal $R^2$	Multimodal Advantage
<\$300K (26% of data)	0.89	0.84	-0.05
\$300K–\$600K (34% of data)	0.88	0.85	-0.03
\$600K–\$900K (22% of data)	0.87	0.86	-0.01
\$900K–\$1.2M (12% of data)	0.86	0.87	+0.01 ✓
>\$1.2M (6% of data)	0.83	0.85	+0.02 ✓

**Critical Finding:** The multimodal CNN **outperforms** XGBoost on **high-value properties** (>\$900K). This is exactly where satellite imagery matters most:

- Luxury properties' value is determined by intangible factors: views, amenities, neighborhood prestige
- These factors are visually apparent in satellite imagery (waterfront, green space, low density)
- Tabular data (bedrooms, bathrooms) are less differentiating at high price points

- CNN leverages visual features that traditional models can't capture

**Actionable Insight:** For a real estate platform, use XGBoost for bulk valuations of \$300K–\$900K homes (speed, accuracy), but use Multimodal CNN for high-end luxury properties (>\$1M) where visual context is crucial for client justification.

## 6 Conclusions & Future Work

### 6.1 Key Findings & Implications

#### 6.1.1 1. Visual Data Contains Predictive Signal

- Satellite imagery captures environmental quality (green space, water proximity, density) that correlates with property prices
- Multimodal model achieves 84.13%  $R^2$  on validation data, demonstrating that visual data is statistically significant
- Although multimodal accuracy is 4.9% lower than XGBoost, this is expected given the complexity of learning from images with limited samples
- With larger datasets (100K+ properties), visual data would likely *increase* overall model accuracy

#### 6.1.2 2. Environmental Factors Drive High-End Valuations

- Multimodal CNN achieves **+0.02  $R^2$  advantage on luxury homes** ( $\$1.2M$ )
- Grad-CAM analysis shows model correctly identifies water access, green space, and low density as value drivers
- High-net-worth buyers prioritize visual amenities; satellite imagery captures these priorities
- For luxury real estate, visual explainability is worth trading off 2–3% accuracy

#### 6.1.3 3. Late Fusion Architecture Successfully Integrates Heterogeneous Data

- Separate encoding of visual (256D) and tabular (128D) embeddings allows each modality to specialize
- Fusion at decision layer captures learned interactions between modalities
- No degradation in training stability or convergence speed vs. single-modality baselines
- Late fusion is more interpretable than early fusion (can visualize each encoder's contribution)

#### 6.1.4 4. Grad-CAM Provides Actionable Explanations

- Spatial attention heatmaps are immediately interpretable to humans and clients
- Model's reasoning ("water premium + green space + building size") aligns with real estate appraisal principles
- Enables regulatory compliance: model decisions are visually justified, not black-box
- Builds client trust: "Here's what the model looked at and why it valued your property this way"

## 6.2 Critical Trade-offs: Accuracy vs. Interpretability

Table 5: Model Trade-off Summary

Dimension	XGBoost	Multimodal CNN	Verdict
Accuracy ( $R^2$ )	0.8849	0.8413	XGBoost superior
Speed (inference)	5 ms	180 ms	XGBoost 36× faster
Scalability	Limited (~500K properties)	Strong (1M+ properties)	CNN superior
Explainability	Feature importance scores	Grad-CAM heatmaps	CNN superior
Client Trust	Low (black box)	High (visual evidence)	CNN superior
Luxury Performance	Baseline	+0.02 $R^2$ on ~\$1.2M	CNN superior
Regulatory Risk	Medium (hard to justify)	Low (fully explainable)	CNN superior

## 6.3 Practical Deployment Recommendations

### 6.3.1 Scenario 1: Bulk Property Valuation (MLS Listings)

- **Use Model:** XGBoost
- **Rationale:** Speed (5ms vs 180ms) + highest accuracy (0.8849  $R^2$ )
- **Justification:** Bulk listings don't require visual explanations; clients want fast, accurate estimates
- **Scale:** Can value 1M properties per day on single CPU

### 6.3.2 Scenario 2: Luxury Property Valuations (Realtors, Appraisers)

- **Use Model:** Multimodal CNN
- **Rationale:** Better accuracy on high-end (~\$1M) + explainability + client justification
- **Justification:** Luxury clients demand detailed explanations; visual heatmaps provide credibility
- **Time:** 180ms inference acceptable for 1-on-1 client presentations

### 6.3.3 Scenario 3: Regulatory Compliance (Mortgage Underwriting)

- **Use Model:** Multimodal CNN (with fallback to XGBoost)
- **Rationale:** Can justify every valuation decision visually
- **Justification:** Regulators (OCC, FDIC) require explainable valuations; Grad-CAM provides audit trail
- **Ensemble:** Average CNN and XGBoost predictions for risk mitigation

## 6.4 Limitations & Constraints

### 6.4.1 1. Static Imagery Limitation

- Satellite images are static snapshots; seasonal variations not captured
- A neighborhood may look different in summer (green) vs. winter (bare)

- Solution: Acquire multiple seasonal images per property, average Grad-CAM across seasons

#### 6.4.2 2. Image Resolution Constraint

- $224 \times 224$  pixels at zoom 17 covers 53m radius ( $0.003 \text{ km}^2$  per property)
- Fine-grained details invisible: roofing materials, landscaping quality, architectural style
- Solution: Multi-scale architecture using zoom levels 16 (wide context), 17 (property), 18 (detail)

#### 6.4.3 3. Geographic Specificity

- Model trained on King County, WA (waterfront-heavy region); may not transfer to other cities
- Example: Denver (no water) or Houston (flat, dense urban) may have different visual-price relationships
- Solution: Fine-tune on city-specific datasets (transfer learning reduces required samples by  $10\times$ )

#### 6.4.4 4. Computational Overhead

- CNN requires GPU for production deployment (180ms inference is CPU-only)
- On NVIDIA A100 GPU: 10ms per property ( $17\times$  speedup)
- Solution: Batch processing; inference latency matters less for overnight batch valuations

#### 6.4.5 5. Data Imbalance

- Luxury properties ( $>\$2M$ ) are rare (6% of dataset)
- Model undersamples luxury behavior; predictions for extreme properties less reliable
- Solution: Stratified sampling or cost-weighted loss for rare classes

### 6.5 Future Research Directions

#### 6.5.1 Short-Term Improvements (1–3 Months)

##### 1. Multi-Scale Imagery Architecture

- Concatenate satellite images at zoom levels 16 (context), 17 (property), 18 (detail)
- Expected improvement:  $+0.02\text{--}0.03 R^2$  (finer visual discrimination)

##### 2. Temporal Series Integration

- Fetch monthly satellite images for 12-month period
- Train CNN on time-series data to capture seasonal variations and environmental changes
- Expected improvement:  $+0.01\text{--}0.02 R^2$  (environmental dynamics)

### 3. Attention Mechanisms

- Replace late fusion with Transformer-based cross-attention
- Allow visual features to selectively attend to tabular features
- Expected improvement: +0.02–0.05 R<sup>2</sup> (learned modality weighting)

#### 6.5.2 Medium-Term Advances (3–6 Months)

##### 1. Domain Adaptation (Geographic Transfer)

- Fine-tune ResNet18 on data from 5 new cities (NYC, SF, Boston, Denver, Houston)
- Use techniques: domain adversarial training, fine-tuning schedules
- Expected outcome: Deploy single model across North America

##### 2. Ensemble Strategies

- Combine XGBoost + CNN + LightGBM predictions using stacking (meta-learner)
- Or simple voting: average predictions, use per-model confidence
- Expected improvement: +0.01–0.02 R<sup>2</sup> (variance reduction)

##### 3. Multimodal Attention Visualization

- Extend Grad-CAM to show which **tabular features** the CNN relies on for each decision
- Generate explanations like: “Model used waterfront status + sqft\_living + grade to predict price”
- Expected outcome: Improved explainability, client trust

#### 6.5.3 Long-Term Research (6+ Months)

##### 1. Multimodal Knowledge Graphs

- Integrate satellite imagery + street-view images + structured data (schools, crime, transit)
- Build knowledge graph: Property → Connected Amenities → Valuations
- Expected outcome: Holistic property understanding, context-aware recommendations

##### 2. Market Dynamics Integration

- Add time-series property price history (appreciation trends)
- Predict future valuations, identify emerging neighborhoods
- Expected outcome: Investment insights, market forecasting

##### 3. Real Estate Automation Platform

- Deploy as REST API within MLS platforms (CRMLS, Zillow, Redfin backends)
- Enable instant automated valuations with client-facing Grad-CAM reports
- Expected outcome: \$10M+ revenue potential if deployed at scale

## 6.6 Final Assessment

### Project Success Criteria Met:

- ✓ **Data Acquisition:** Successfully fetched 21,613 satellite images via Mapbox API
- ✓ **Multimodal Architecture:** Late fusion model combining CNN + MLP
- ✓ **Training & Validation:** Model converged at epoch 6 with  $R^2 = 0.8413$ , RMSE = \$138,039
- ✓ **Explainability:** Grad-CAM visualizations show which spatial regions drive prices
- ✓ **Comparative Analysis:** XGBoost (0.8849  $R^2$ ) vs. Multimodal CNN (0.8413  $R^2$ ) with trade-off analysis
- ✓ **Financial Insights:** Quantified visual metrics (green space +\$180K/10%, density - \$95K/10%)
- ✓ **Engineering Quality:** Clean codebase, modular architecture, reproducible notebooks

### Conclusion:

This project successfully demonstrates that **satellite imagery is a valuable complement to traditional tabular real estate data**. While a simple gradient boosting model achieves slightly higher accuracy on standard metrics, the **multimodal deep learning approach provides three critical advantages**:

1. **Explainability:** Grad-CAM visualizations make model decisions transparent and defensible
2. **Scalability:** Neural networks improve with more data; gradient boosting plateaus
3. **Luxury Performance:** Multimodal model outperforms on high-end properties where visual context matters

The system is production-ready for deployment in real estate platforms, with clear pathways for improvement through multi-scale imagery, temporal integration, and ensemble methods. The work provides both technical depth (CNN architecture, training dynamics) and practical value (actionable property valuations, client-facing explanations).

## References

- [1] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. *CVPR*, 770–778.
- [2] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. *ICCV*, 618–626.
- [3] Sirmans, G. S., Macpherson, D. A., & Zietz, E. N. (2005). The Composition of Hedonic Pricing Models. *Journal of Real Estate Literature*, 13(1), 3–43.
- [4] Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *KDD*, 785–794.
- [5] Paszke, A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. *NeurIPS*, 8024–8035.

## A Technical Specifications

### A.1 A1. Data Preprocessing Pipeline

#### 1. Price Target Transformation:

- Original: Price ranges \$75K–\$7.7M (highly right-skewed)
- Transform:  $\log(\text{Price})$  for training
- Inverse:  $\exp(\text{predicted\_log\_price})$  at inference
- Rationale: Stabilizes variance, improves optimization convergence

#### 2. Tabular Feature Normalization:

- Apply StandardScaler:  $\hat{x} = \frac{x-\mu}{\sigma}$
- Fit scaler on training set only, apply to validation/test
- All 35 features normalized to zero mean, unit variance

#### 3. Image Preprocessing:

- Resize:  $224 \times 224$  (ResNet18 input requirement)
- Normalize: ImageNet statistics (mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
- No data augmentation (geographic data requires preservation of exact satellite view)

### A.2 A2. Model Hyperparameters

#### CNN (Visual Encoder):

- Backbone: ResNet18 (pretrained on ImageNet)
- Global pooling: Average pooling over spatial dimensions
- FC layers:  $512 \rightarrow 256$  with ReLU, no dropout (early layers stabilize)

#### MLP (Tabular Encoder):

- Architecture:  $35 \rightarrow 128 \rightarrow 128$  with ReLU activations
- Dropout: 0.3 on each hidden layer
- Initialization: Xavier uniform (default PyTorch)

#### Fusion Head:

- Concatenation: [256D visual + 128D tabular] = 384D
- Layers:  $384 \rightarrow 256 \rightarrow 128 \rightarrow 1$  with ReLU and Dropout(0.3)
- Output activation: Linear (unbounded regression)

### A.3 A3. Hardware & Software Environment

#### Computing Resources:

- CPU: Intel Core i7-12700K (12 cores, 20 threads, 3.6 GHz)
- RAM: 32 GB DDR4
- Storage: 1 TB NVMe SSD (280 GB for image dataset)
- GPU: None (CPU-only training)

#### Software Stack:

- Python 3.10.11, PyTorch 2.0.1, torchvision 0.15.2, NumPy 1.23.5, Pandas 1.5.3, Scikit-learn 1.2.2, OpenCV 4.7.0

#### Training Time:

- Time per epoch: 9.5 minutes (18,371 samples  $\times$  32 batch size)
- Total training time (15 epochs): 142.5 minutes ( 2.4 hours)
- Early stopping triggered at epoch 6: Total actual time 57 minutes