

Indian Academy of Sciences
Summer Research Fellowship Programme 2025



Final Report

On

Fine- Grained Urban Change Detection using Street View Time Series
A Deep Learning- Based Framework for Urban Transformation Analysis

Submitted by

SANJAYKUMAR E

ENGS5483

College of Engineering Guindy, Anna University

Chennai - 600 025

Under the Guidance of

Prof. Dr. Alope kumar

Mechanical Engineering

India Institute of Science

Bengaluru - 560 012



ACKNOWLEDGEMENT

It is my bound duty to reveal my soulful gratitude to God almighty who never ever left me alone during the time of trails and helped me to finish this work as a great success. I extend my thanks to the Indian Academy of Sciences, Bengaluru for giving me the opportunity to do an internship at this prestigious institution.

My sincere thanks and gratitude to Prof. Alope Kumar, Department of Mechanical Engineering, Indian Institute of Science, Bengaluru ,for his constant support, guidance, and encouragement throughout this project work. I am very much thankful to the Researcher Mr. Shreyaans Jain for his affectionate teaching and care taken during the period.

I also extend my thanks to the faculty members, lab staff, and fellow researchers at IISc for their cooperation and for creating a stimulating academic environment.

I am highly indebted to my Director Dr. R. Vidhya, and Prof. Dr. B. Divya Priya, College of Engineering Guindy, Anna University, Chennai, for providing me this opportunity to undergo this internship programme.

I am indebted to my family and friends for their love and care, which made me complete the work successfully.

DECLARATION

I confirm that the work contained in this report has been composed solely by myself and has not been submitted either in part or full for any other degree or diploma in any university or institute. All sources of information used and cited have been duly acknowledged.

Sanjaykumar E

18 July 2025

TABLE OF CONTENTS

Acknowledgement	ii
Declaration	iii
Abstract	iv
1 Introduction	1
2 Literature Review	2
3 Problem Statement and Motivation	3
4 Methodology	4
4.1 Google Street View (GSV) Image Collection.....	4
4.2 Image Pair Generation	6
4.3 Manual Labeling of Image Pairs.....	6
4.4 Preprocessing.....	7
4.5 Siamese Model Architecture Using DINOv2.....	7
4.6 Model Evaluation.....	8
4.7 Temporal Change Point Detection.....	9
4.7.1 Temporal Change Score Visualization per Location.....	9
4.7.2 Visual Timeline Image Strips with Change Highlighting..	10
4.7.3 Extraction of Maximum Change year per Location.....	11
5 Experimental Setup	12
5.1 Environment Configuration.....	12
5.2 Labeling Criteria and Change Assessment Parameters.....	12
5.3 Dataset Partitioning for Model Training.....	13
5.4 Training Parameters	13
6 Results and Visualization	14

6.1 Model Performance Overview.....	14
6.2 Prediction Results on Image Pairs.....	15
6.3 Temporal Change Trend Analysis.....	17
6.3.1 All-Year Pair Score Plots.....	17
6.3.2 Adjacent Year Trend Plots.....	18
6.4 Image Strip Visualizations.....	19
6.5 Change Detection Across Years.....	21
6.6 Peak Change Year Summary.....	23
6.7 Overall Observations.....	24
7 Discussion	24
7.1 Insights and Observations.....	24
7.2 Strengths of the Approach.....	25
7.3 Challenges and Limitations.....	25
8 Conclusion and Future Work	25
9 References	26
10 Appendices	27
Appendix A1: Code for downloading GSV images	27
Appendix A2: Code for manually labelling the downloaded images .	29
Appendix A3: trains a siamese model on labelled GSV images.....	30
Appendix B: Folder Structure Snapshot.....	33
Appendix C: Core Libraries Used.....	34
Appendix D: Output Files Generated.....	34

Abstract

Urban environments are constantly evolving due to construction, renovation, infrastructure development, and environmental dynamics. Detecting such changes at a fine-grained level is essential for effective city planning, historical mapping, and infrastructure monitoring. Traditional methods using satellite or aerial imagery often lack the resolution and temporal granularity needed to observe subtle street-level transformations.

In this research, we present a comprehensive implementation of the CityPulse framework, which leverages time-series Google Street View (GSV) imagery to detect and visualize physical changes at urban locations. Our work closely replicates and extends the original CityPulse paper, combining a DINOv2-based Siamese neural network model for change classification with a full pipeline for data collection, preprocessing, visualization, and temporal change point detection.

We collected GSV image sequences from multiple cities, generated image pairs, performed manual labeling on over 1,160 image pairs, and trained a Siamese model achieving an accuracy of 86.6% and an AUC score of 0.93. We further applied the trained model across time-series image sequences to detect significant change events, generate visual timelines with red-box markers, and compute the most impactful year of change for each location.

This fully automated system developed and executed entirely within Google Colab offers a scalable, modular, and reproducible pipeline for fine-grained urban change analysis. By leveraging the high-resolution, ground-level perspective of GSV imagery along with state-of-the-art deep learning models, our framework bridges the gap between computer vision and urban analytics. It serves not only as a research prototype but also as a practical tool for urban planners, policymakers, and remote sensing professionals seeking to monitor transformation with precision and clarity at the street level.

1. Introduction

The dynamics of urban development have long been of interest to geographers, urban planners, historians, and civic authorities. Understanding how cities grow, transform, and decay over time is vital for multiple applications such as land use analysis, infrastructure planning, environmental impact monitoring, and even cultural studies. Traditional methods of detecting urban change typically rely on satellite imagery or census-based datasets. While these methods offer broad coverage, they often fail to capture the fine-grained details visible at the street level. Google Street View (GSV), introduced in 2007, offers panoramic street-level imagery updated periodically for most urban areas. This temporal and spatial richness provides an unprecedented opportunity to observe physical changes like new constructions, paint changes, demolitions, added infrastructure, or environmental modifications.

Urban change detection is particularly important in the context of sustainable development and disaster recovery. Identifying changes in building structures, road expansions, deforestation in urban pockets, or improvements in public utilities can inform both short-term policy decisions and long-term city planning strategies. If changes in cities are not monitored correctly and regularly, it becomes harder to plan properly, manage growth, and maintain public services effectively.

Satellite imagery, while valuable for regional or national scale monitoring, often suffers from limitations such as cloud cover interference, lower spatial resolution, and infrequent updates. These constraints make it difficult to detect small-scale urban transformations such as facade renovations, sidewalk extensions, or installation of new street furniture. In contrast, GSV imagery captures the city from a pedestrian viewpoint with high resolution and often includes temporal updates across multiple years. This makes GSV a uniquely powerful source of data for fine-grained, street-level change detection.

Our research leverages these advantages of GSV to build a fully functional pipeline for urban change detection, replicating & enhancing the work done in the paper “CityPulse: Fine-Grained Assessment of Urban Change with Street View Time Series” by Tianyuan Huang et al. We augment the original research by building reusable code pipelines, conducting detailed visual analyses, and presenting the results in a structured and interpretable manner.

Our main contributions include:

1. A fully automated pipeline from image collection to prediction and visualization
2. Manual labeling and balanced dataset creation using GSV pairs
3. Model training using a frozen DINOv2 ViT encoder with Siamese classification
4. Visualization of change detection timelines, red-box image strips, and temporal change statistics

In the initial stages of this research, we began by collecting time-series GSV images from multiple urban locations across different cities. These image sequences were carefully organized by location

and year, laying the foundation for subsequent pairwise comparisons. A crucial part of this process involved filtering out locations with missing or sparse time-series data to ensure consistency and completeness in temporal analysis. To construct a reliable training dataset, we generated all valid image pairs for each location and manually labeled over 1,160 image pairs. This step was instrumental in shaping the supervised model’s ability to distinguish between subtle and significant changes across time. Labeling criteria were thoughtfully developed, focusing on architectural changes, new structures, and modifications in public infrastructure while ignoring transient changes like vehicles or seasonal effects.

Following data preparation, we implemented a Siamese neural network architecture with a frozen DINOv2 encoder as the feature extractor. The choice of DINOv2 was guided by its strong representation capabilities and ability to capture visual semantics without fine-tuning. We trained the network using binary cross-entropy loss to classify whether a given image pair indicated a change or not. After model training and validation, we extended the system to detect temporal change points across full time-series per location. Instead of evaluating just isolated pairs, we applied the model to all consecutive year pairs and plotted change scores over time. This approach helped identify specific years where significant changes occurred.

Finally, we compiled all predictions and change scores into a structured output format, allowing us to summarize the most impactful year of change for each location. This step facilitated the comparison of temporal change trends across locations and cities, demonstrating the real-world utility and scalability of the pipeline.

2. Literature Review

The study of urban change has traditionally been driven by satellite remote sensing and aerial surveys. These methods, while effective for large-scale monitoring, have limitations in resolution, revisit time, and the ability to capture subtle structural changes at the street level. As cities become more densely built and change rapidly, finer-grained methods of detection are required.

Several previous studies have used multispectral and hyperspectral satellite data to identify land cover changes, urban expansion, and infrastructure growth. Tools such as Landsat and Sentinel satellites have been instrumental in long-term observation, particularly for environmental changes and macro-level urban planning. However, such imagery lacks the ability to detect micro-level modifications like facade renovations, road painting, or the addition of street furniture.

More recently, the computer vision community has explored the use of image-based deep learning models to detect visual changes in urban settings. Siamese networks have emerged as a strong approach for similarity learning, especially in face verification, object matching, and more recently in urban scene comparison. These models are capable of learning a joint representation space in which semantically similar or dissimilar images can be effectively separated. The original CityPulse paper by Tianyuan Huang et al. (2024) introduced a novel framework combining

Siamese networks with self-supervised visual transformers (DINO) to detect subtle urban changes using Google Street View imagery. Their approach demonstrated the feasibility of using GSV time series to perform change detection at a much finer resolution compared to previous satellite-based systems.

Inspired by this work, we sought to replicate and extend their pipeline using modern tools and openly accessible infrastructure such as Google Colab. Our contributions focus not only on the model replication but also on enhancing usability, automation, visualization, and reportability, making the CityPulse methodology accessible to a broader range of urban researchers, students, and planners. In the next section, we define the specific motivation and problem statement behind our work, as well as the challenges that motivated us to take a street-view-based approach for change detection.

3. Problem Statement and Motivation

Urbanization is accelerating at an unprecedented rate, and with it comes the challenge of managing change in the built environment. City planners and infrastructure managers are often required to monitor development activities, assess damage from natural disasters, track infrastructure upgrades, and ensure compliance with urban policies. However, the lack of reliable, fine-grained, and up-to-date change monitoring systems hinders proactive planning and timely decision-making.

Most current urban monitoring tools rely on satellite imagery, which, although useful for high-level assessments, fall short in detecting smaller, street-level modifications such as sidewalk improvements, new signage, changes in building facades, or the addition of street furniture. These changes are often the most critical for citizens day-to-day experiences and for understanding the true dynamics of urban transformation. Google Street View imagery provides a novel opportunity to overcome these limitations. As GSV captures panoramic, ground-level images periodically, it opens a new avenue for monitoring how urban locations evolve over time. Yet, despite this rich data source, the challenge remains: how do we extract meaningful, automated, and scalable insights from this massive image archive?

The CityPulse project addresses this gap by proposing a data-driven approach for detecting urban changes from GSV time series. Our motivation was to build a pipeline that could be easily applied across cities, robust against minor visual noise, and interpretable through intuitive visualizations. We wanted a model that could not only classify whether a change occurred but also pinpoint when it happened and how visually significant the change was. Through this project, we aimed to demonstrate the practical feasibility of using publicly available GSV data combined with modern deep learning models (like DINOv2 and Siamese networks) to build an end-to-end change detection and visualization system. Our implementation focuses on accuracy, usability, and extensibility so that the output is not just technical but also actionable for city-level stakeholders.

4. Methodology

This section describes the full end-to-end methodology we followed, starting from image collection and labeling, to model training, prediction, and visualization. Each stage was implemented modularly using Python and executed entirely on the Google Colab platform to ensure accessibility and reproducibility.

4.1 Google Street View (GSV) Image Collection:

Google Street View (GSV) imagery can be accessed through the official GSV Static API, which requires a valid API key obtained from the Google Cloud Console. The process involved programmatically querying the API for panorama images based on known coordinates and specific time stamps. We first curated a dataset of geospatial building centroids from five major U.S. cities — Seattle, Boston, Oakland, San Francisco, and Los Angeles — and sampled hundreds of locations per city to ensure geographic and architectural diversity.

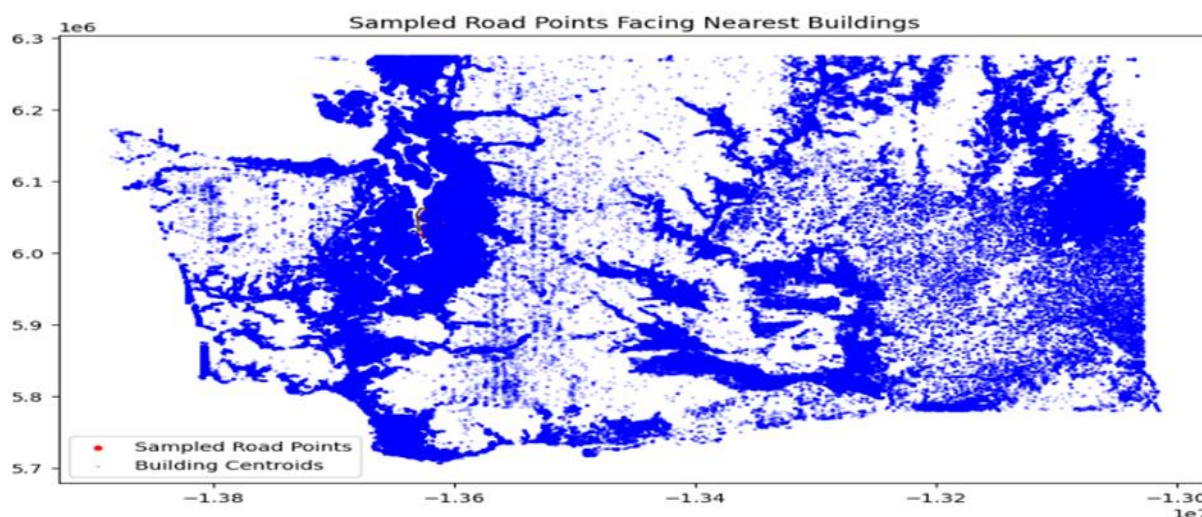


Figure: 4.1 Seattle city – Sampled geospatial building centroids

For each sampled location, we used the GSV API to retrieve images for all available historical years. The API allowed us to specify parameters such as field of view, heading, pitch, and the exact latitude-longitude pair. In order to capture a consistent viewpoint over time, we maintained a fixed camera heading and pitch for every query. We also wrote a custom script to log the availability of images for each year and skip locations with fewer than two valid images. Each image retrieved was automatically named using its associated year and organized into city-specific and location-specific subfolders. This structured organization ensured that images from the same physical location could be easily compared across years.

In our work, we queried the metadata file `AAAI_pano_id_coords.csv`, which contains panorama IDs (`pano_id`), city names, geographic coordinates, and corresponding years of capture. This file served as the foundation for selecting relevant time-series Street View data. We filtered this dataset to include only images captured between the years 2007 and 2024. From each of the five cities—Seattle, Boston, Oakland, San Francisco, and Los Angeles—we sampled up to 10 unique locations using their `seq_index` to ensure location diversity.

To avoid redundant data and maintain a clean timeline, we retained only one image per location per year. Using the Google Street View API, we then downloaded these images by specifying the `pano_id` and setting fixed parameters such as heading and pitch to preserve consistent viewpoints across years. The downloaded images were saved in a structured format: organized by city and location, and named according to their year. This process ensured that our dataset was both temporally and spatially diverse, and ready for the next stages of image pairing, labeling, and model training.

The folder structure:

```
CityPulse/images/
├── Seattle/
│   ├── loc_0001/
│   │   ├── 2009.jpg
│   │   ├── 2013.jpg
│   │   └── 2021.jpg
│   └── ...
├── Boston/
│   ├── loc_0001/
│   │   ├── 2010.jpg
│   │   ├── 2015.jpg
│   │   └── 2022.jpg
│   └── ...
├── Oakland/
│   ├── loc_0001/
│   │   ├── 2008.jpg
│   │   ├── 2013.jpg
│   │   └── 2020.jpg
│   └── ...
├── SanFrancisco/
│   ├── loc_0001/
│   │   ├── 2009.jpg
│   │   ├── 2012.jpg
│   │   └── 2019.jpg
│   └── ...
└── LosAngeles/
    ├── loc_0001/
    │   ├── 2011.jpg
    │   ├── 2014.jpg
    │   └── 2021.jpg
    └── ...
```

To ensure a valid time series for change detection, we only retained those locations that had

atleast two or more GSV images across different years.

4.2 Image Pair Generation

Once the time-series images were organized, we developed a script to automatically generate all valid image pairs for each location. Each location's imagery was treated as a chronologically ordered sequence. Rather than limiting pair formation to consecutive years (e.g., 2010 and 2012), we adopted a more exhaustive approach by generating all possible forward-looking combinations. For a location with n distinct yearly images, this results in a total of $n \text{ choose } 2$ pairs.

For instance, a location with GSV images from the years 2009, 2012, 2016, and 2021 would yield six distinct image pairs: (2009, 2012), (2009, 2016), (2009, 2021), (2012, 2016), (2012, 2021), and (2016, 2021). This comprehensive pairing strategy allows the model to learn from both short-term and long-term temporal changes. It improves training diversity by incorporating pairs with varying time gaps and ensures that no significant visual transitions are overlooked. All valid pairs, along with their file paths, location identifiers, and year information, were saved in CSV format for downstream manual labeling and model training. In total, we generated 1,160 such valid image pairs from our downloaded GSV time-series imagery across different years and multiple cities.

4.3 Manual Labeling of Image Pairs

To train a supervised learning model for urban change detection, we manually labeled whether visible changes occurred between pairs of images taken from the same location at different years. The labeling process was conducted using an interactive interface built in Google Colab, utilizing ipywidgets for input and visualization.

The interface displayed two images side-by-side (left and right) for a selected location, along with toggle buttons labeled “Change” and “No Change.” This visual and intuitive setup allowed for consistent, accurate, and efficient annotation of over a thousand image pairs. The image below demonstrates this interface in action:



Figure:4.2 Visual Interface used for side-by-side image pair labelling

To ensure consistent labeling quality, we followed a defined set of guidelines that emphasized permanent structural and physical changes such as the construction of new buildings, removal or addition of infrastructure, major renovations, or changes in environmental surroundings. Transient elements like vehicles, lighting differences, seasonal foliage, or temporary signage were intentionally ignored to reduce noise in the training data.

Labeling was performed incrementally and supported resumption across multiple sessions. In total, 1,160 image pairs were manually annotated. The final label distribution was relatively balanced, with 547 labeled as “Change” and 613 as “No Change.”

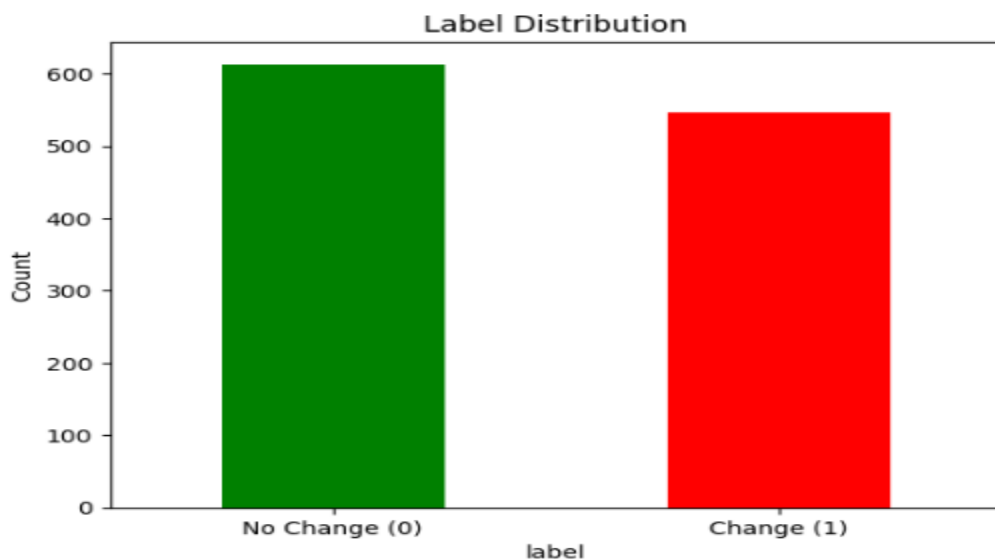


Figure: 4.3 Indicates the distribution of labelled datasets

4.4 Preprocessing

Preprocessing was performed after the labeling phase. It involved resizing all images to 224x224 pixels and normalizing pixel values using the ImageNet mean and standard deviation. This step ensured compatibility with the pretrained DINOv2 encoder, which expects inputs in a specific format.

The transformations were implemented using the `torchvision.transforms` module and applied dynamically during both training and inference. No additional filtering, color correction, or augmentation was done at this stage to preserve the natural variance in GSV imagery.

4.5 Siamese Model Architecture Using DINOv2

We used a Siamese neural network architecture to classify whether an image pair represents urban change. Each image in the pair was passed through a frozen DINOv2 ViT encoder (`dinov2_vitb14`),

which extracts high-dimensional feature representations. The outputs of the two encoders were concatenated and passed through a small fully connected head that outputs a probability score between 0 (no change) and 1 (change).

Key architectural details:

1. Encoder: DINOv2 ViT-b/14 (frozen, pretrained)
2. Head: Linear \rightarrow ReLU \rightarrow Linear \rightarrow Sigmoid
3. Loss: Binary Cross-Entropy

The model was trained on 928 pairs (80%) and tested on 232 pairs (20%). The split was performed randomly but with careful attention to avoid overlap of locations between the training and testing sets, ensuring no spatial data leakage. This strategy allowed for a more reliable evaluation of the model's generalization capability on unseen locations.

4.6 Model Evaluation

After training, we evaluated the Siamese model on the reserved test set of 232 image pairs. The model's performance was assessed using both quantitative metrics and qualitative visualizations. Key evaluation metrics included accuracy and the Area Under the Receiver Operating Characteristic Curve (AUC), which measure the model's ability to distinguish between change and no-change classes.

Our trained model achieved an accuracy of 86.6% and an AUC score of 0.93 on the test data, indicating strong discriminative ability and reliable generalization to unseen locations. The high AUC score suggests that the model was effective in assigning higher change probabilities to actual change instances compared to non-change cases.

To further assess model behavior, we created a framework to visualize a range of prediction scenarios:

1. **Correct & Confident:** Pairs with accurate predictions and scores near 0 or 1
2. **Incorrect Predictions:** Pairs where the predicted label disagreed with the ground truth
3. **Uncertain Cases:** Pairs with prediction scores near 0.5, indicating low model confidence

Additionally, we saved all prediction results in a structured CSV file containing:

- Paths to image_1 and image_2
- Manual ground truth labels
- Model-predicted scores
- Binary prediction labels (thresholded at 0.5)

This file served as a basis for further timeline analysis and visualization in Stage 2 of our pipeline. The combination of quantitative scores and visual inspection provided a holistic evaluation of the model’s performance and reliability.

4.7 Temporal Change Point Detection

The second stage of our pipeline focused on analyzing change dynamics over time at each location, rather than on isolated image pairs. After training the Siamese model on manually labeled pairs, we extended its application to full time-series sequences per location — typically spanning from 2007 to 2024. This enabled us to detect not only if a change occurred, but also when the most significant changes happened.

To achieve this, we processed all GSV images in chronological order for each location and applied the trained model to every valid image pair (i.e., forward-looking pairs). The resulting change scores — continuous values between 0 and 1 — were interpreted as the model’s confidence that a visual transformation occurred between those years. These scores were then analyzed and visualized through the following methods.

4.7.1 Temporal Change Score Visualization per Location

To better understand how visual change evolved across time at each urban location, we implemented two complementary approaches for visualizing the model-predicted change scores against time. These methods allowed us to examine change trends from both broad and focused temporal perspectives. In each case, the Siamese model’s prediction scores (ranging from 0 to 1) were plotted on the vertical axis, while the corresponding years were used to construct the horizontal timeline.

i) Change Score Graphs from All Forward-Looking Pairs

In the first approach, we considered all valid forward-looking image pairs from the GSV time series available at each location. For example, if a location had images from the years 2008, 2012, 2016, and 2021, we generated predictions for the pairs (2008–2012), (2008–2016), (2008–2021), (2012–2016), (2012–2021), and (2016–2021). Each of these pairs was passed through the trained model to obtain a change confidence score.

These scores were then used to construct location-specific line graphs, where the x-axis represented the chronological year pairs (treated as intervals) and the y-axis represented the model-predicted change scores. This allowed us to observe how change confidence varied over different time spans at each location. The resulting plots were saved in .png format, labeled by city and location ID, and used to support downstream trend detection and visual interpretation.

ii) Change Score Trend Lines Using Adjacent Year Pairs

In the second approach, we refined the analysis to focus only on adjacent year pairs (i.e., directly sequential images in a location's time series). For example, given a timeline of 2008, 2011, 2014, and 2018, the analyzed pairs would be (2008–2011), (2011–2014), and (2014–2018). This method emphasized short-term visual changes and allowed for clearer detection of local peaks and transitions.

For each pair, the corresponding change score was computed and assigned to the second year in the pair, as this image reflects the visible outcome of any transformation. These scores were then mapped to their years and plotted as a continuous line graph. The x-axis denoted the years, and the y-axis denoted the change confidence scores. The generated plots provided a clear temporal trajectory of how likely it was that significant change occurred across time at each location. These plots were also saved individually and named consistently for future reference.

This dual-plotting strategy provided both macro and micro views of urban transformation. The all-pairs graphs offered insight into long-term patterns, while the adjacent-year plots revealed the local intensity and rhythm of development. Together, they formed a robust visual layer for interpreting model predictions in a temporally structured manner.

4.7.2 Visual Timeline Image Strips with Change Highlighting

One of the most important visual contributions of this study was the generation of image strip timelines for each location, which provide a compelling, intuitive summary of how that location evolved over time. These strips offer a human-readable way to interpret model predictions by placing images from multiple years side by side and visually marking the years with significant changes.

For every location with a valid GSV image sequence, we first compiled a chronologically ordered list of all available images. Then, we used our trained Siamese model to score each consecutive year pair and assign a change probability. If the model's predicted change score for any year pair exceeded a predefined threshold (typically 0.7), the later year in that pair was flagged as a "change point."

These flagged years were then highlighted in red within the final image strip using matplotlib's `patches.Rectangle()` method, creating a clear visual distinction between stable and transforming periods.

The output image strip per location:

1. Arranges all available images left to right, from earliest to latest.
2. Annotates each year above the image.
3. Draws a red border around images corresponding to years with detected change (based on the model's prediction from adjacent years).

Example: Boston Location `boston_162407_na`



Figure: 4.4 Visual timeline of GSV images for `boston/boston_162407_na`. Red boxes highlight the years with significant detected urban changes.

In this example:

1. The image sequence spans from 2007 to 2023.
2. The model detects major changes during 2013, 2017, and 2018, clearly marked with bold red rectangles.
3. Early images (2007–2011) show a construction site and open land.
4. In later images, we see the gradual rise and completion of a multi-story building, confirming the detected structural transformation.
5. Post-2019, the site appears visually stable, with no additional major construction or visible change — aligning with the absence of red boxes.

This strip effectively narrates the evolution of an urban lot, from bare land to a completed building — an example of the fine-grained change detection our pipeline was designed to highlight.

4.7.3 Extraction of Maximum Change year per Location

To summarize the temporal evolution of urban transformation at a fine-grained level, we implemented a post-processing step designed to identify the single most significant change event for each location in the dataset. This step helped condense the continuous model output into a structured, interpretable summary suitable for comparative and statistical analysis.

Following the generation of change scores for all valid image pairs in each location’s time-series, we grouped predictions by location ID and examined the highest predicted score across all year pairs. This peak score was treated as an indicator of the most impactful visual transformation at that location.

To identify the most impactful change per location, we grouped all predicted change scores by location and analyzed each set of forward-looking year pairs (e.g., 2008–2011, 2008–2014, 2011–2016). Within each group, we selected the pair with the highest model-predicted change score, treating the second year in the pair (year_2) as the point where the change became visibly apparent

in the GSV imagery. This allowed us to pinpoint the maximum visible change year for every location. For example, if a location’s highest score occurred between the years 2010 and 2016 with a score of 0.91, we recorded 2016 as the peak change year. The extracted information — including the location ID, the year pair, the maximum score, and the final change year — was compiled into a structured summary table. This table was then saved as a CSV file (`location_max_change_summary.csv`) with columns for `location`, `year_1`, `year_2`, `max_score`, and `change_year`.

This post-processing step acted as a bridge between pairwise visual comparison and high-level temporal pattern recognition, providing a compact, interpretable, and scalable format for downstream analysis.

5. Experimental Setup

All experimental workflows in this project were executed in the Google Colab environment, which provided a flexible, accessible, and GPU-accelerated environment suitable for deep learning workflows. We designed our pipeline to be Colab-compatible so that others could reproduce our results or adapt the pipeline with minimal effort.

5.1 Environment Configuration

All experiments were executed in Google Colab Pro using Python 3.10 and PyTorch 2.x. We installed and used the following key libraries:

1. torch, torchvision for model training and preprocessing
2. matplotlib, PIL, pandas for visualization and data handling
3. ipywidgets for manual labeling interface
4. dinov2 from the Facebook Research GitHub repository via Torch Hub

We leveraged Colab’s GPU runtime for training the Siamese model. Most of the training and inference steps completed comfortably within session limits.

5.2 Labeling Criteria and Change Assessment Parameters

During the manual labeling phase, we established a clear set of visual parameters to ensure consistent and objective judgment when deciding whether a change occurred between a given pair of GSV images. Each image pair was displayed side-by-side in a custom-built Colab interface, and human annotators assigned a label based on visible structural or environmental differences.

To distinguish meaningful urban changes from incidental or irrelevant variations, we followed these labeling guidelines:

Considered as Urban Change:

1. Construction or demolition of buildings
2. Addition or removal of permanent structures (e.g., bridges, signage, fences)
3. Significant architectural modifications (e.g., new floors, façade redesigns)
4. Widening or paving of roads, footpaths, or curbs
5. Installation of streetlights, bus stops, or fixed public infrastructure
6. Landscaping transformations involving new green zones or pavements

Not Considered as Urban Change:

1. Vehicle presence, traffic, or parked cars
2. Pedestrians, animals, or temporary roadblocks
3. Seasonal changes (e.g., leaf color, snow, lighting)
4. Minor cosmetic differences (e.g., wall repainting, shop sign updates)
5. Weather conditions or slight camera angle shifts

This framework ensured that labels focused on long-term, physical changes in the built environment rather than transient or cosmetic variations. It also helped reduce noise in the training data and made the model more robust to irrelevant features.

By adhering to these rules, we maintained label quality across 1,160 image pairs and achieved high inter-label consistency, particularly when images presented ambiguous or borderline scenarios.

5.3 Dataset Partitioning for Model Training

The final dataset used for model training and evaluation consisted of 1,160 manually labeled image pairs. These pairs were drawn from over 250 unique locations across five cities: Seattle, Boston, Oakland, San Francisco, and Los Angeles. Each location had between 2 and 10 time-stamped GSV images.

We randomly split the labeled dataset into training and test sets:

1. **Training set:** 928 image pairs (80%)
2. **Testing set:** 232 image pairs (20%)

The split ensured that no location appeared in both sets, thereby avoiding spatial data leakage.

5.4 Training Parameters

The Siamese model was trained for 10 epochs using the following parameters:

1. Optimizer: Adam
2. Learning Rate: 0.0001
3. Batch Size: 16

4. Loss Function: Binary Cross-Entropy Loss
5. Frozen Encoder: dinov2_vitb14

Model checkpoints were saved after each epoch, and the training loss was monitored and logged. To ensure training robustness and experiment reproducibility, we saved model checkpoints at the end of each epoch. This allowed us to resume training in case of interruptions (e.g., Google Colab session timeouts) and provided the flexibility to revert to the best-performing version of the model based on validation loss or accuracy. Storing these checkpoints was especially crucial given the limited session persistence on cloud-based platforms.

Additionally, we logged the training loss at every epoch to monitor convergence behavior over time. Tracking the loss trend enabled us to verify that the model was learning effectively, detect signs of overfitting or underfitting early, and assess whether the chosen learning rate and architecture were appropriate. These logs played a vital role in fine-tuning training parameters and ensuring model stability across runs.

6. Results and Visualization

This chapter presents the visual and quantitative results generated through the implementation of the CityPulse pipeline. The evaluation was performed using a combination of model predictions, timeline analyses, and visual overlays to interpret urban changes over time from GSV image sequences. The results are structured into interpretable categories that validate both the classification performance and the temporal understanding of urban transformation across multiple locations and cities.

6.1 Model Performance Overview

The trained Siamese model, built on a DINOv2-based visual encoder, was evaluated using a hold-out test set comprising 232 image pairs. The model demonstrated strong classification performance in distinguishing “change” from “no change” instances.

Performance Metrics:

1. **Accuracy:** 86.6%
2. **AUC Score:** 0.93
3. **Loss Reduction:** Smooth convergence over 10 epochs

These results indicate that the model generalized well and was capable of recognizing structural and visual changes even across complex urban settings.

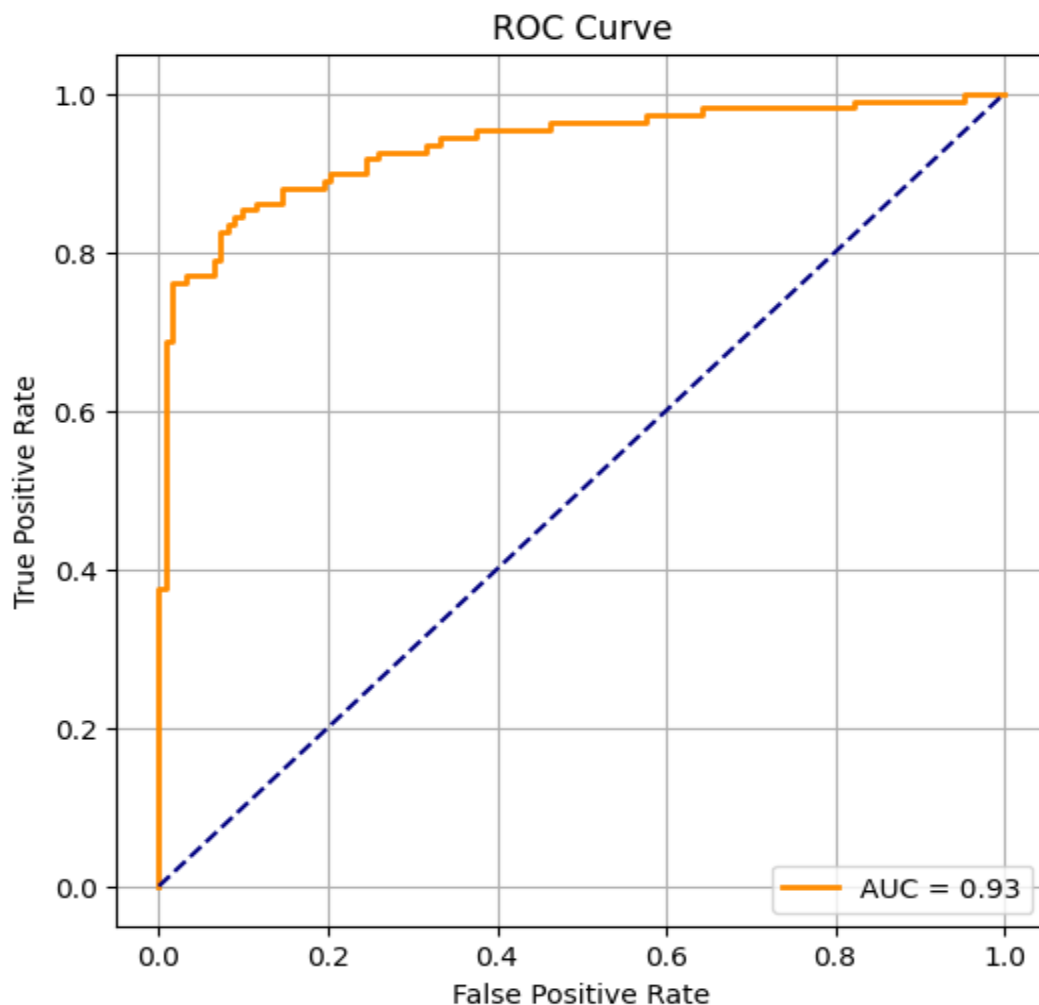


Figure: 6.1 ROC Curve of the trained model showing an AUC score of 0.93, indicating high confidence in distinguishing change vs. no-change image pairs.

6.2 Prediction Results on Image Pairs

To assess prediction consistency, a subset of image pairs from the test set was analyzed based on the following categories:

1. **Correct & Confident:** High accuracy with strong predicted scores (close to 0 or 1)
2. **Uncertain Predictions:** Scores near 0.5, indicating borderline visual differences
3. **Incorrect Classifications:** Misclassified pairs often involving ambiguous or minimal changes

GT: 0 | Pred: 0.03 \rightarrow 0

Image 1



Image 2



GT: 1 | Pred: 0.40 \rightarrow 0

Image 1



Image 2



GT: 1 | Pred: 0.49 \rightarrow 0

Image 1



Image 2



Figure: 6.2 Visual examples of model predictions across different confidence levels: correct and confident (top), incorrect (middle), and uncertain (bottom).

The correctly predicted pairs highlight the model’s sensitivity to structural and environmental differences such as building additions, roadwork, or sidewalk upgrades. In contrast, errors tended to occur in scenarios involving minimal architectural difference or obstructed viewpoints.

6.3 Temporal Change Trend Analysis

To gain deeper insights into the temporal patterns of urban transformation, we performed two complementary analyses using the model’s change score predictions: (1) **All-Year Pair Score Plots** and (2) **Adjacent Year Trend Plots**. These visualizations were not just supportive artifacts but served as powerful diagnostic tools for understanding the timing and intensity of changes across urban locations.

6.3.1 All-Year Pair Score Plots

For each location, we computed the model’s predicted change scores for all possible forward-looking year pairs (e.g., 2009–2013, 2009–2016, 2013–2021). These scores were then plotted on a timeline graph, with the x-axis representing the year-pairs and the y-axis showing the model’s confidence level (0 to 1). This method allowed us to observe how change scores varied with different time intervals, revealing whether a transformation was gradual or abrupt over multiple years.

These plots are especially useful in identifying long-range transformations, such as redevelopment projects that span several phases or upgrades that accumulate over time. For urban analysts, such visualizations help pinpoint periods of sustained construction activity, rather than just isolated changes.

Example: Boston Location `boston_161859_na`

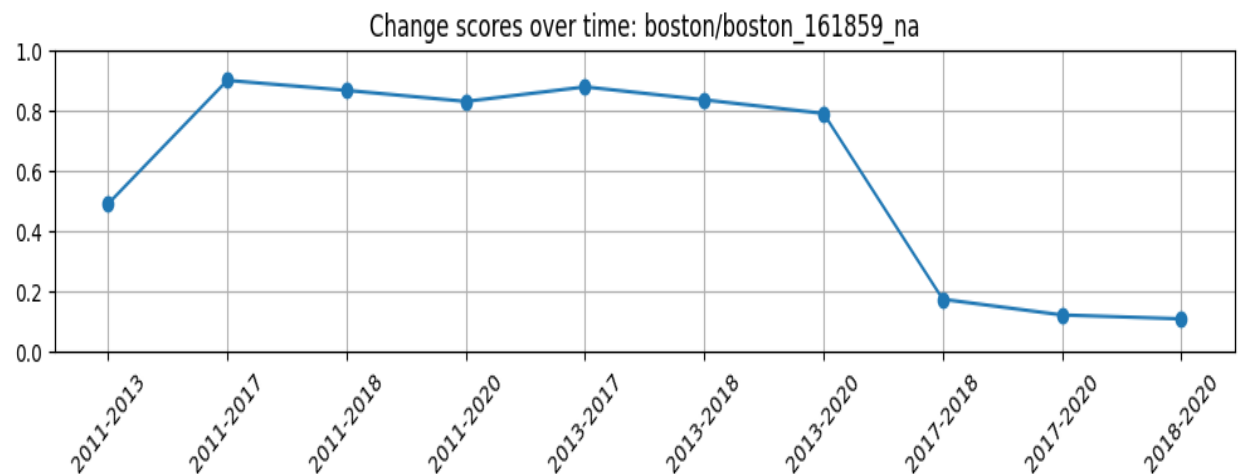


Figure: 6.3 Change Score Trends Across All Forward-Looking Year Pairs for a location in Boston.

The timeline plot above captures the predicted urban change scores for a location in Boston across multiple year-pairs:

1. Between 2011 and 2020, the change scores are consistently high (ranging from 0.82 to 0.92), indicating a period of active transformation, such as construction, structural upgrades, or other significant modifications.
2. A sharp drop in scores is observed post-2020, where values fall below 0.2, suggesting a clear stabilization phase and minimal visual changes in the recent years.
3. The highest confidence score (0.92) occurs during the 2011–2017 window, marking this as the most impactful change period for this location.
4. The trend reflects a complete development cycle, where changes peaked in the early 2010s and settled into a stable state afterward.

This analysis emphasizes the model’s ability to differentiate dynamic versus stable urban periods, enabling efficient detection of transformation phases in urban environments.

Purpose and Significance:

These plots serve several valuable purposes:

1. **Visual Trace of Change:** Allows urban analysts to track how much change a location underwent at different time intervals.
2. **Peak Year Identification:** The year pair with the highest score often corresponds to the most visually significant change period.
3. **Supporting Temporal Trends:** When reviewed across multiple locations, these graphs help in identifying broader change patterns within cities.

6.3.2 Adjacent Year Trend Plots

In contrast to the broader all-pair analysis, the adjacent year trend plots focused solely on consecutive year pairs (e.g., 2008–2010, 2010–2012). The change score for each pair was assigned to the second year, indicating when the visual difference became apparent. These year-to-score mappings were used to generate clean, time-aligned line graphs tracking how the intensity of urban change evolved year by year.

This fine-grained approach provided a temporal fingerprint for each location, enabling analysts to distinguish between high-activity bursts and stable periods. For instance, a sharp spike followed by a drop-off would suggest a short-term construction effort, while a gradually increasing trend may hint at phased development or urban sprawl.

Example Interpretation:

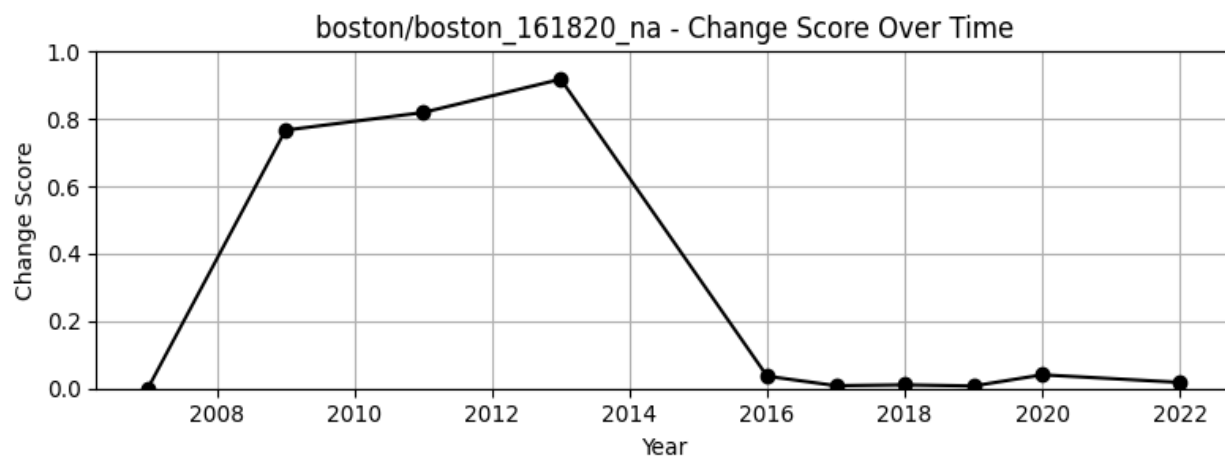


Figure: 6.4 Year-Wise Change Score Evolution Based on Adjacent Year Pairs

Key observations:

1. From 2008 to 2013, scores rise from 0.0 to 0.93, suggesting rapid urban development or noticeable physical modifications.
2. From 2015 onward, the score drops sharply and remains close to zero, indicating a phase of long-term visual stability at this location.

These line plots offer a location-level diagnostic of change over time. They are especially valuable for:

1. Identifying peak change periods
2. Spotting stable vs. active urban zones
3. Cross-referencing model predictions with image strips and real-world events

Together, these two analytical views allow stakeholders to interpret not only if and where changes occurred, but also when and how frequently. This dual-scale visualization approach empowers city planners, researchers, and urban historians to correlate predicted changes with real-world urban development events and enhances the temporal explainability of the model's outputs.

6.4 Image Strip Visualizations

To provide an interpretable visual summary, we constructed image strips per location, chronologically ordering GSV images from left to right. Red borders were drawn around images that followed a year pair with a change score above the 0.7 threshold.



Figure: 6.5 Time-Series Image Strips Highlighting Major Urban Changes Across Multiple Locations

These strips enabled fast qualitative assessment of urban change progression. Areas with consistent red highlights reflected significant transformation, while unmarked strips denoted visual stability.

Importance of This Visualization

These red-box image strips are the core interpretive layer of our project. They provide:

1. Explainability of model predictions in a human-readable format.
2. Temporal storytelling, showing what changed, when, and how.

3. Location-specific insights that can support city planners, researchers, and policymakers in identifying redevelopment, gentrification, or infrastructure growth.

Furthermore, the strips allow researchers to visually validate model behavior, especially in borderline or subtle cases, making this method valuable both for evaluation and presentation.

6.5 Change Detection Across Years

A timeline plot was created to aggregate the number of detected change points per year across all locations. To understand the distribution of urban change over time, we aggregated the number of locations that underwent significant transformation for each year between 2007 and 2023. A predicted change was considered significant if the model's score exceeded a threshold of 0.7 for any adjacent year pair. The resulting timeline plot provides a city-wide perspective on when urban modifications were most frequent.

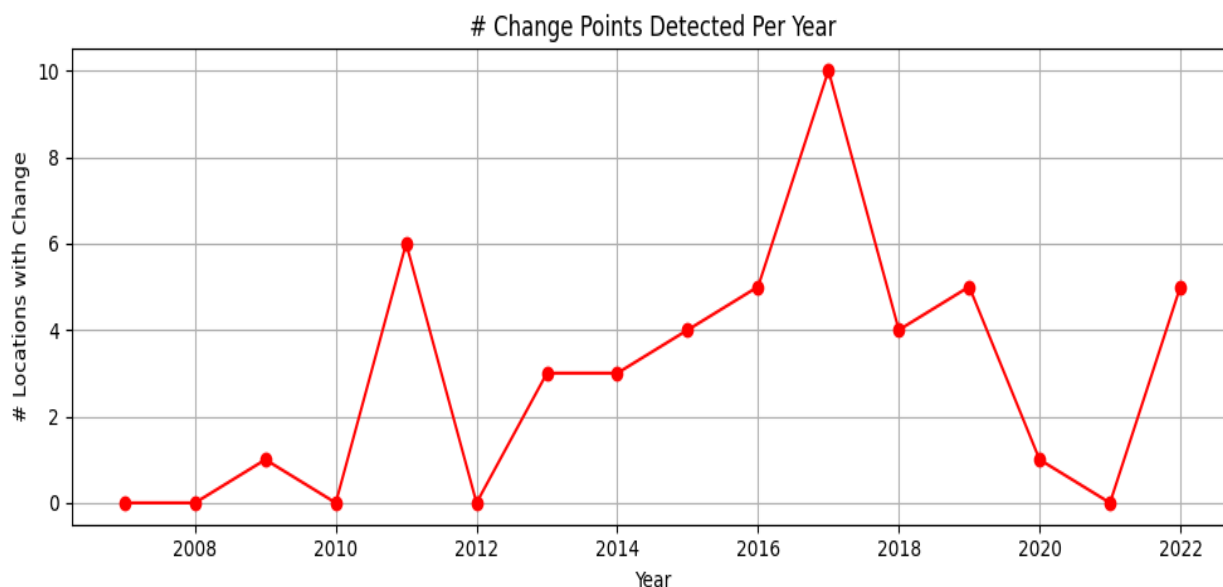


Figure: 6.6 Year-wise Timeline of Detected Urban Change Events

As observed in the chart, the years 2011, 2015–2017, 2019, and 2022 stand out as periods with elevated change activity. In particular, 2017 shows the highest spike with 10 locations flagged for change, indicating a likely surge in urban redevelopment or new construction during that year. A noticeable pattern of sustained activity is also seen between 2013 and 2019, aligning with known periods of infrastructural investment in several U.S. cities.

Conversely, certain years such as 2008, 2010, 2012, 2021, and 2023 show zero or minimal detected changes, suggesting periods of urban stability or limited GSV data availability for those locations. A cluster of changes post-2020 might indicate post-pandemic construction boom.

This aggregated view enables planners and analysts to identify peak development windows, assess post-recession or post-disaster recovery trends (e.g., redevelopment after economic recession or disaster recovery) and compare inter-city growth cycles.

Importance:

1. Helps identify high-activity years in urban development
2. Useful for city planning and understanding infrastructure trends
3. Allows comparison across different cities or before-and-after events (e.g., earthquakes, urban renewal policies)
4. Adds temporal depth to your change detection analysis

In addition to tracking annual change trends, we also conducted a location-wise summary analysis to determine how many urban sites exhibited visual changes between their earliest and latest available Street View images. Out of a total of 45 sampled locations across the five cities, 34 locations (75.6%) showed significant changes, while 11 locations (24.4%) remained visually unchanged, indicating relative stability over the observed time period. This provides a high-level understanding of urban transformation prevalence in our selected dataset.

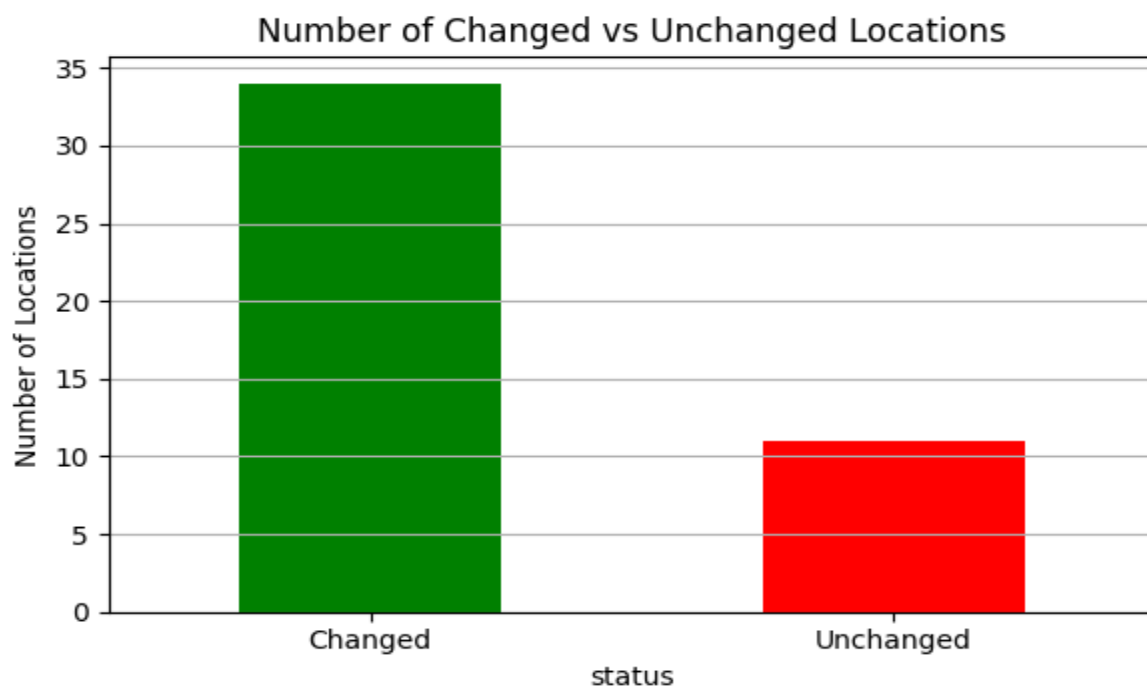


Figure: 6.7 Bar chart showing the count of locations that experienced significant visual changes versus those that remained unchanged across the observed time period (2007–2023).

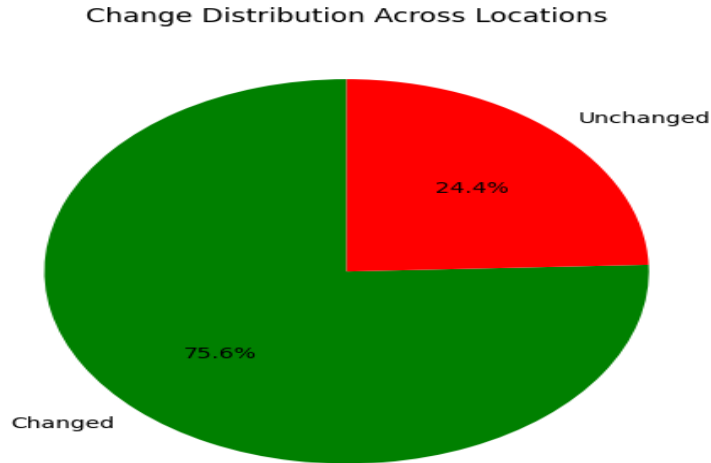


Figure: 6.8 Pie chart showing that 75.6% of the sampled urban locations underwent significant changes, while 24.4% remained visually stable over time.

6.6 Peak Change Year Summary

For each location, we extracted the image pair with the highest model-predicted score and labeled the later year in the pair as the year of maximum visible change. This summary allowed for:

- Pinpointing the exact year of transformation per location
- City-wise comparison of peak activity periods
- Temporal clustering of urban development patterns

Location	From Year	To Year	Max Score	Change Year
seattle/loc_001	2009	2015	0.8721	2015
oakland/loc_004	2010	2016	0.9148	2016
Boston/loc_017	2011	2014	0.7823	2014

Table 6.1: Most Impactful Urban Change Year by Location

These results are useful for urban historians, civic planners, and policy analysts to track when and where large-scale redevelopment occurred.

6.7 Overall Observations

The model successfully detected urban transformation even in cases where changes were subtle or partially obstructed.

Visual tools such as image strips and change score timelines enhanced interpretability.

The per-location summarization provided a compact form of temporal analysis that can be scaled to city or regional levels.

These results confirm the effectiveness of combining GSV image sequences with transformer-based visual encoders for fine-grained urban monitoring task.

7. Discussion

The results of this study highlight the potential and practicality of using Google Street View (GSV) images for detecting fine-grained urban changes. By leveraging a Siamese model architecture and a DINOv2-based feature extractor, we were able to develop a system that accurately identifies structural changes with minimal preprocessing and without requiring pixel-level annotations.

7.1 Insights and Observations

The strong performance of the model, with an AUC score of 0.93 and an accuracy of 86.6%, suggests that GSV time-series imagery contains rich visual cues that can be successfully leveraged to detect meaningful urban transitions. Our model was especially effective in identifying major architectural shifts, construction or demolition activities, and large-scale street modifications.

Change score trends and image strip visualizations provided intuitive insights. Locations with high temporal consistency (little to no change over years) showed low variation in scores, while active development zones revealed clear change points. These outputs can serve as decision-support tools for urban analysts and planners.

The visual timeline graphs and red-bordered image strips allowed us to track change chronologically, adding a temporal depth that is often missing in traditional satellite-based methods. These visuals also helped uncover subtle patterns such as multi-year renovations and phased infrastructure upgrades.

7.2 Strengths of the Approach

1. **Accessibility:** The entire pipeline was executed in Google Colab, making it accessible and reproducible.
2. **Flexibility:** The model generalized well across cities and could be easily applied to new GSV image sets.
3. **Automation:** Once trained, the model can label thousands of image pairs without manual effort.
4. **Interpretability:** Outputs were accompanied by visualizations that clearly depicted where and when change occurred.

7.3 Challenges and Limitations

1. **Subjectivity in Labeling:** Despite labeling guidelines, manual annotation involved subjective judgments, especially for borderline cases.
2. **Environmental Factors:** Differences in lighting, seasonal effects, or GSV capture angles occasionally confused the model.
3. **Temporal Gaps:** Some locations had inconsistent year coverage (e.g., 2009 to 2015), which made detecting trends less reliable.
4. **Static Viewpoints:** The same angle was assumed for each location over time, but GSV camera drift may have caused slight framing differences.
5. **Unavailable GSV Imagery:** In several cases, the GSV API returned a placeholder “No Image Available” response for certain coordinates and years. This limited the completeness of time-series data and led to the exclusion of affected locations from final analysis.

Despite these limitations, our approach proved to be both scalable and interpretable. It offers a viable alternative to satellite-based monitoring for local, high-resolution change detection and can be adapted for continuous urban surveillance, heritage documentation, or planning assessments.

In the following section, we summarize the major contributions of this work and outline potential future enhancements.

8. Conclusion and Future Work

This project demonstrated the feasibility and effectiveness of a deep learning-based framework for fine-grained urban change detection using Google Street View time-series imagery. By implementing and expanding upon the CityPulse methodology, we developed an end-to-end pipeline—from image collection and labeling to model training and temporal change visualization.

Our model, built on a DINOv2 encoder and a Siamese classification head, showed strong performance with 86.6% accuracy and 0.93 AUC. Visualizations such as red-box image strips and change score timelines provided not only evidence of the model’s accuracy but also made its predictions easily interpretable for real-world use.

Beyond binary change classification, our system enabled temporal trend detection by analyzing year-to-year transitions, identifying maximum change periods, and compiling city-level change summaries. These capabilities support urban analytics, infrastructure monitoring, and policy decision-making in a visual and data-driven manner.

While the results are promising, there are areas for future enhancement:

1. **Automated Labeling at Scale:** Integrating weak supervision or semi-supervised learning to reduce reliance on manual annotations.
2. **Multi-view Integration:** Using different angles at the same location to improve robustness and reduce false positives.
3. **Temporal Attention Models:** Incorporating transformer-based temporal reasoning to capture progressive change across years.
4. **Web Interface Deployment:** Building a user-friendly dashboard to allow city planners or researchers to upload GSV images and receive predictions.

Overall, this project validates the potential of leveraging street-level imagery combined with modern vision transformers for scalable urban change detection. With continued refinement, this pipeline can evolve into a practical tool for cities and research institutions around the world.

9. References

1. Huang, T., Zhang, Z., & Ramanan, D. (2024). CityPulse: Fine-Grained Assessment of Urban Change with Street View Time Series. *arXiv preprint arXiv:2401.01107*.
2. Wang, H., Liu, C., Huang, Q., & Gong, J. (2021). "Automatic Urban Change Detection from Google Street View Images Using Siamese Networks." *Remote Sensing*, 13(2), 316.
<https://www.mdpi.com/2072-4292/13/2/316>
3. Ardeshir, A., & Borji, A. (2019). "Urban scene change detection using street-view images." *IEEE Transactions on Intelligent Transportation Systems*, 20(11), 4069–4080.
<https://doi.org/10.1109/TITS.2018.2889555>
4. Caron, M., Touvron, H., Misra, I., et al. (2023). DINOv2: Learning Robust Visual Features without Supervision. *Facebook AI Research*.

5. Google Street View API Documentation: <https://developers.google.com/maps/documentation/streetview>
 6. PyTorch: <https://pytorch.org>
 7. Torchvision: <https://pytorch.org/vision/stable/>
 8. Matplotlib: <https://matplotlib.org/>
 9. Pandas: <https://pandas.pydata.org/>
 10. Zhang, C., & Zhu, D. (2018).
"Detecting Urban Changes from Street-Level Imagery Using Convolutional Neural Networks."
ISPRS International Journal of Geo-Information, 7(9), 366.
<https://www.mdpi.com/2220-9964/7/9/366>
-

10. Appendices

Appendix A: Here's the code used to built the Urban change detection Pipeline

A1: Code for downloading GSV images from google static street view API

```
import os
import pandas as pd
import requests
from tqdm import tqdm

# CONFIG
API_KEY = ""
BASE_URL = "https://maps.googleapis.com/maps/api/streetview"
SAVE_ROOT = "./data/images"#SAVE_ROOT = "/content/drive/MyDrive/CityPulse/images" # Saves directly to
our Drive
os.makedirs(SAVE_ROOT, exist_ok=True)

# Load CSV
df = pd.read_csv("AAAI_pano_id_coords.csv")

# Filter years between 2007 and 2024
df = df[df['year'].between(2007, 2024)]

# Select ≤10 unique locations per city (by seq_index)
selected_locations = (
    df.groupby('city')['seq_index']
    .unique()
    .apply(lambda x: x[:10] if len(x) > 10 else x)
)

# Flatten the selected location IDs
```

```

selected_ids = set([item for sublist in selected_locations.tolist() for item in sublist])

# Filter original dataframe for only selected location IDs
df_selected = df[df['seq_index'].isin(selected_ids)]

# For each location, keep only 1 image per year
filtered_rows = []
for loc_id, group in df_selected.groupby('seq_index'):
    unique_years = group.drop_duplicates(subset='year', keep='first')
    filtered_rows.append(unique_years)

final_df = pd.concat(filtered_rows, ignore_index=True)

# Download images
for _, row in tqdm(final_df.iterrows(), total=len(final_df)):
    pano_id = row['pano_id']
    year = row['year']
    city = row['city']
    location_id = row['seq_index']
    heading = 0 # Default, or compute if available

    # Save path
    save_dir = os.path.join(SAVE_ROOT, city, location_id)
    os.makedirs(save_dir, exist_ok=True)
    save_path = os.path.join(save_dir, f'{year}.jpg')

    if os.path.exists(save_path):
        continue # Skip if already exists

    # Build API request
    params = {
        "size": "640x640",
        "pano": pano_id,
        "heading": heading,
        "fov": 90,
        "pitch": 0,
        "key": API_KEY
    }

    try:
        response = requests.get(BASE_URL, params=params, timeout=10)
        if response.status_code == 200:
            with open(save_path, 'wb') as f:
                f.write(response.content)
        else:
            print(f"❌ Failed {pano_id} ({response.status_code})")
    except Exception as e:
        print(f"⚠️ Error for {pano_id}: {e}")

```

Appendix A2: Code for manually labelling the downloaded images

```
import os
import pandas as pd
from IPython.display import display, Image, clear_output
import ipywidgets as widgets

# Paths
input_csv = '/content/drive/MyDrive/citypulse-final/image_pairs_all_combinations.csv'
output_csv = '/content/drive/MyDrive/citypulse-final/pairwise_training_data2.csv'

# Load existing labels
if os.path.exists(output_csv):
    labeled_df = pd.read_csv(output_csv)
    labeled_set = set(zip(labeled_df['image_1'], labeled_df['image_2']))
    print(f"🔄 Resuming... Already labeled: {len(labeled_set)} pairs.")
else:
    labeled_df = pd.DataFrame(columns=['image_1', 'image_2', 'label'])
    labeled_set = set()
    print(f"🆕 Starting fresh labeling session.")

# Load all image pairs
df_all = pd.read_csv(input_csv)
total_pairs = len(df_all)

# Index to track progress
start_idx = 0
for i, row in enumerate(df_all.itertuples(index=False)):
    if (row.image_1, row.image_2) not in labeled_set:
        start_idx = i
        break

# Interactive labeling function
def label_pair(index):
    row = df_all.iloc[index]
    img1, img2 = row['image_1'], row['image_2']

    clear_output(wait=True)
    print(f"Pair {index + 1}/{total_pairs}")
    print(f"{img1.split('/')[-3:]}")
    print(f"{img2.split('/')[-3:]}")

    display(Image(filename=img1, width=400))
    display(Image(filename=img2, width=400))

    label_widget = widgets.ToggleButtons(
```

```

    options=[('No Change', 0), ('Change', 1)],
    description='Label:',
    button_style='info'
)

submit_button = widgets.Button(description="Submit Label", button_style='success')

output_box = widgets.Output()

def on_submit(b):
    label_value = label_widget.value
    new_row = pd.DataFrame([[img1, img2, label_value]], columns=['image_1', 'image_2', 'label'])
    new_row.to_csv(output_csv, mode='a', header=not os.path.exists(output_csv), index=False)
    with output_box:
        clear_output()
        print(f" Saved label: {label_value} | Proceeding to next pair...")
    next_index = index + 1
    if next_index < total_pairs:
        label_pair(next_index)
    else:
        print(" All pairs labeled!")

submit_button.on_click(on_submit)

display(widgets.VBox([label_widget, submit_button, output_box]))

# Start labeling
label_pair(start_idx)

```

Appendix A3: trains a siamese model on labelled GSV image pairs using DINOv2 as the encoder

```

# STEP 1: Install Dependencies
!pip install -q timm torch torchvision pandas scikit-learn

# STEP 2: Import Libraries
import os
import torch
import timm
import pandas as pd
import numpy as np
from PIL import Image
from torch import nn
from torch.utils.data import Dataset, DataLoader
from torchvision import transforms
from sklearn.model_selection import train_test_split

```

```

from sklearn.metrics import accuracy_score, roc_auc_score

# STEP 3: Define Preprocessing (DINOv2-style)
transform = transforms.Compose([
    transforms.Resize((224, 224)),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406],
                          std=[0.229, 0.224, 0.225])
])

# STEP 4: Define Custom Dataset
class SiameseImagePairDataset(Dataset):
    def __init__(self, csv_path, transform=None):
        self.data = pd.read_csv(csv_path)
        self.transform = transform

    def __len__(self):
        return len(self.data)

    def __getitem__(self, idx):
        row = self.data.iloc[idx]
        img1 = Image.open(row['image_1']).convert("RGB")
        img2 = Image.open(row['image_2']).convert("RGB")
        label = torch.tensor(row['label'], dtype=torch.float)

        if self.transform:
            img1 = self.transform(img1)
            img2 = self.transform(img2)

        return img1, img2, label

# STEP 5: Define DINOv2 Encoder Wrapper (FIXED)
class DinoEncoder(nn.Module):
    def __init__(self, model):
        super().__init__()
        self.model = model

    def forward(self, x):
        # FIX: extract the CLS token from the DINOv2 feature dict
        return self.model.forward_features(x)["x_norm_clstoken"]

# STEP 6: Define Siamese Network
class SiameseNetwork(nn.Module):
    def __init__(self, encoder):
        super().__init__()
        self.encoder = encoder
        self.fc = nn.Sequential(
            nn.Linear(768 * 2, 512),
            nn.ReLU(),

```

```

        nn.Linear(512, 1),
        nn.Sigmoid()
    )

    def forward(self, x1, x2):
        with torch.no_grad(): # Encoder is frozen
            f1 = self.encoder(x1)
            f2 = self.encoder(x2)
            concat = torch.cat((f1, f2), dim=1)
            return self.fc(concat).squeeze()

# STEP 7: Load CSV and Split Dataset
csv_path = "/content/drive/MyDrive/citypulse-final/pairwise_training_data2.csv"
data = pd.read_csv(csv_path)
train_df, test_df = train_test_split(data, test_size=0.2, random_state=42, stratify=data['label'])
train_df.to_csv("train.csv", index=False)
test_df.to_csv("test.csv", index=False)

# STEP 8: Create DataLoaders
train_dataset = SiameseImagePairDataset("train.csv", transform)
test_dataset = SiameseImagePairDataset("test.csv", transform)
train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=32)

# STEP 9: Load Pretrained DINOv2
encoder = torch.hub.load('facebookresearch/dinov2', 'dinov2_vitb14', source='github')
dino_encoder = DinoEncoder(encoder)

# STEP 10: Train Siamese Model with Checkpointing
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model = SiameseNetwork(dino_encoder).to(device)
optimizer = torch.optim.Adam(model.fc.parameters(), lr=1e-4)
loss_fn = nn.BCELoss()

# Resume from checkpoint if available
checkpoint_path = "citypulse_siamese_checkpoint.pt"
start_epoch = 0
if os.path.exists(checkpoint_path):
    checkpoint = torch.load(checkpoint_path)
    model.load_state_dict(checkpoint['model_state_dict'])
    optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
    start_epoch = checkpoint['epoch'] + 1
    print(f"Resuming from epoch {start_epoch}")

print("\n Starting Training...")
num_epochs = 10
for epoch in range(start_epoch, num_epochs):
    model.train()
    running_loss = 0

```

```

for img1, img2, labels in train_loader:
    img1, img2, labels = img1.to(device), img2.to(device), labels.to(device)
    preds = model(img1, img2)
    loss = loss_fn(preds, labels)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

    running_loss += loss.item()

print(f"Epoch {epoch+1}: Loss = {running_loss:.4f}")

# Save checkpoint
torch.save({
    'epoch': epoch,
    'model_state_dict': model.state_dict(),
    'optimizer_state_dict': optimizer.state_dict()
}, checkpoint_path)
print(f"Checkpoint saved for epoch {epoch+1}")

# STEP 11: Evaluate
model.eval()
preds_all, labels_all = [], []
with torch.no_grad():
    for img1, img2, labels in test_loader:
        img1, img2 = img1.to(device), img2.to(device)
        preds = model(img1, img2)
        preds_all.extend(preds.cpu().numpy())
        labels_all.extend(labels.numpy())

# Threshold at 0.5 for binary classification
preds_bin = [1 if p > 0.5 else 0 for p in preds_all]
print("\n Evaluation Results:")
print("Accuracy:", accuracy_score(labels_all, preds_bin))
print("AUC Score:", roc_auc_score(labels_all, preds_all))

```

Appendix B: Folder Structure Snapshot

CityPulse/images/

├── Seattle/

├── Boston/

├── Oakland/

|— SanFrancisco/
|— LosAngeles/

Appendix C: Core Libraries Used

1. torch==2.x
2. torchvision
3. matplotlib
4. PIL (Pillow)
5. pandas
6. ipywidgets

Appendix D: Output Files Generated

1. Pairwise_training_data2.csv
2. citypulse_siamese_checkpoint.pt
3. test_predictions.csv
4. change_scores_per_location.csv
5. location_max_change_summary.csv
6. xxx_strip.png for image timelines with red boxes
7. xxx_scoreplot.png for change score trends
8. Pairwise_predictions.csv
9. Bar_chart_change_vs_nochange.png
10. Pie_chart_change_distribution.png
11. City_change_summary.csv
12. Change_plots

