

Structure from Motion System in Modern Browser

BY XU PU

HOME PAGE - *ptx.digital*

EMAIL - *ptx.pluto@gmail.com*

Abstract

This thesis aims at designing a 3D reconstruction system on the modern web platform, including a structure from motion pipeline and its integration/optimization in the web platform.

1 Introduction

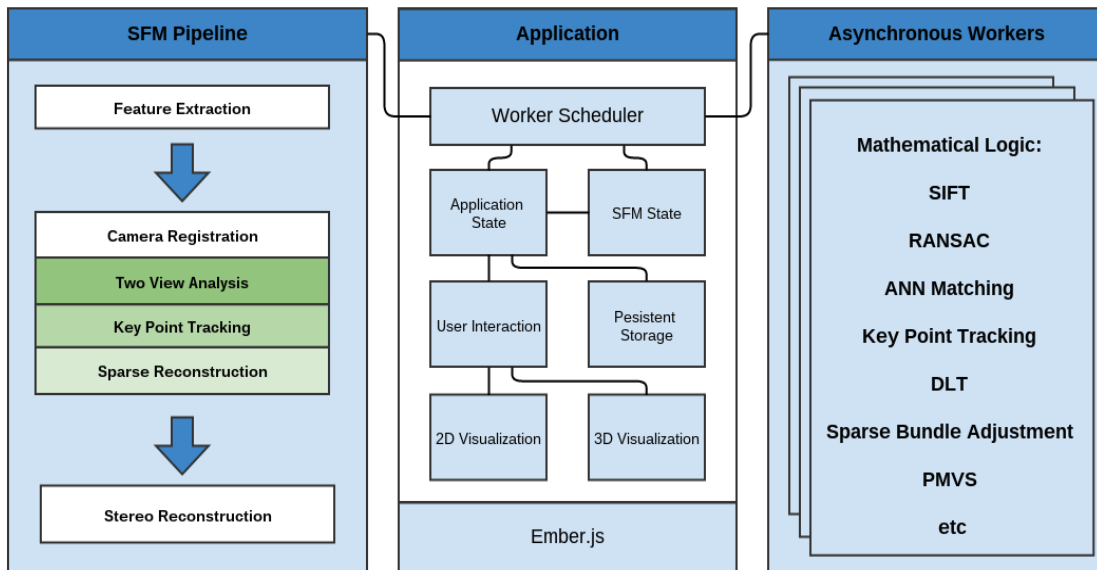
The title *STRUCTURE FROM MOTION SYSTEM IN MODERN BROWSER* implies two key aspects, *3D Reconstruction* and *Web Platform*, and both sides have great potential.

On one hand, structure from motion has always been a critical part of computer vision. In recent years, structure from motion has become increasingly relevant because of the rise of consumer drones and VR. Reconstruct 3D scene directly from photos taken by regular digital camera instead of specialized equipments can dramatically decrease the requirements for this important task, and open up the door to potentially unlimited resource for 3D reconstruction.

On the other hand, modern browser as a platform is gaining its popularity. It has become increasingly powerful over the years both in functionality and performance. Several novel features of modern browser has made implementing an SFM system possible. First is *WebWorkers*, javascript is a single threaded language, *WebWorkers* enables parallel computing in client-side javascript. Another one is *WebGL*, which is an interface for *OpenGL*. It made rendering of 3D models inside browser possible. Indexed DB is also a very important element, it implements a light weight object storage database, enables versatile persistent storage inside browser. Last but not the least, *ArrayBuffer* and *TypedArray* made handling large dataset using javascript practical.

This project attempts to construct a 3D reconstruction system on the modern web platform, take advantage of the advanced capability of modern browser and its cross platform nature, make 3D reconstruction available everywhere.

2 Overview



. Framework of WebSFM

3 Geometry Context

3.1 Coordinate System

Points in the pipeline have multiple coordinates according to different frames of reference.

- Image Coordinate – x
- Camera Coordinate – X_{cam}
- World Coordinate – X

Coordinate transitions:

- Rotation Matrix(R) and Translation Vector(t) – $X_{\text{cam}} = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix} \cdot X$
- Calibration Matrix(K) – $x = K \cdot [I, 0] \cdot X_{\text{cam}}$
- Projection Matrix(P) – $x = PX = K \cdot [R, t] \cdot X$

3.2 Camera Model

Camera model have multiple forms

- Standard Linear Projective Camera Model (Undistorted)
 - Rotation Matrix (R)
 - Translation Vector (t)
 - Calibration Matrix (K)
 - Projection Matrix (P=K[R,t])
- Parameterized Full Camera (with distortion) - 11 variables
 - Euler Angles (r) - 3 variables
 - Translation Vector (t) - 3 variables
 - Principal Point (px, py) - 2 variables
 - Focal Length (f) - 1 variable
 - Radial Distortion (k1, k2) - 2 variables

Parameterized full camera with distortion is only used during camera registration. Distortion is ignored before camera registration; after registration cameras are undistorted.

4 Pipeline

Structure from motion pipeline of WebSFM is largely based on Bundler.

First generate SIFT features for each image (details in 4.1), then perform robust two view analysis on each pair of images (details in 4.2), matches produced by this process can then be used to track key points.

A *track* represents a point in the 3D space, it has one or no view on an image, we will find tracks by merging two view matches to a graph and traversing the match graph. Tracks that have more than one view on a single image is excluded as an obvious outlier.

Tracks are then used to initiate camera registration. We first choose a pair of camera to initiate using *five point algorithm*[4] to estimate first pair of camera parameter, then triangulate tracks they observe and run a two frame *sparse bundle adjustment* (details in 4.4) to refine the initial estimation. After initiation, incremental recovery is performed (details in 4.3), and sparse reconstruction is obtained.

Currently WebSFM does not handle reconstruction beyond sparse. PMVS/CMVS[6] can generate dense oriented point cloud, which can then be used to recover surface polygon using *Poisson Surface Reconstruction*[2].

4.1 SIFT (Scale Invariant Feature Transform)

Lowe's SIFT[1] is our choice of feature extractor. SIFT features can be detected and matched regardless of its scale, orientation or lighting. Following is an overview of SIFT algorithm:

Algorithm SIFT

```

features := empty array
for each octave:
  scales space := progressive blur the base image of current octave
  dog space := subtract adjacent scales in scale space
  for each extrema in dog space:
    if principal curvature of extrema is too big or negative: break
    subpixel := interpolate extrema position in dog space, at (x,y,scale)
    orientations := dominant orientations of gradient histogram in scale space
    for each orientation in orientations:
      descriptor := gradient distribution histogram in scale space
      features += (octave,subpixel,orientation,descriptor)
    end for
  end for
end for

```

Octave space start at -1 or 0, choose -1 if image is too small. In each octave, we will scan 3 DoG layers, which requires 5 layers in DoG space, and 6 layers in Gaussian space. When blurring the image, larger sigma means larger kernel size, then directly leads to more time consumption, so blur is done progressively to keep sigma relatively small. We use default descriptor size stated in lowe's original paper, $4 \times 4 \times 8$ histogram and 128 dimensional Uint8 vector.

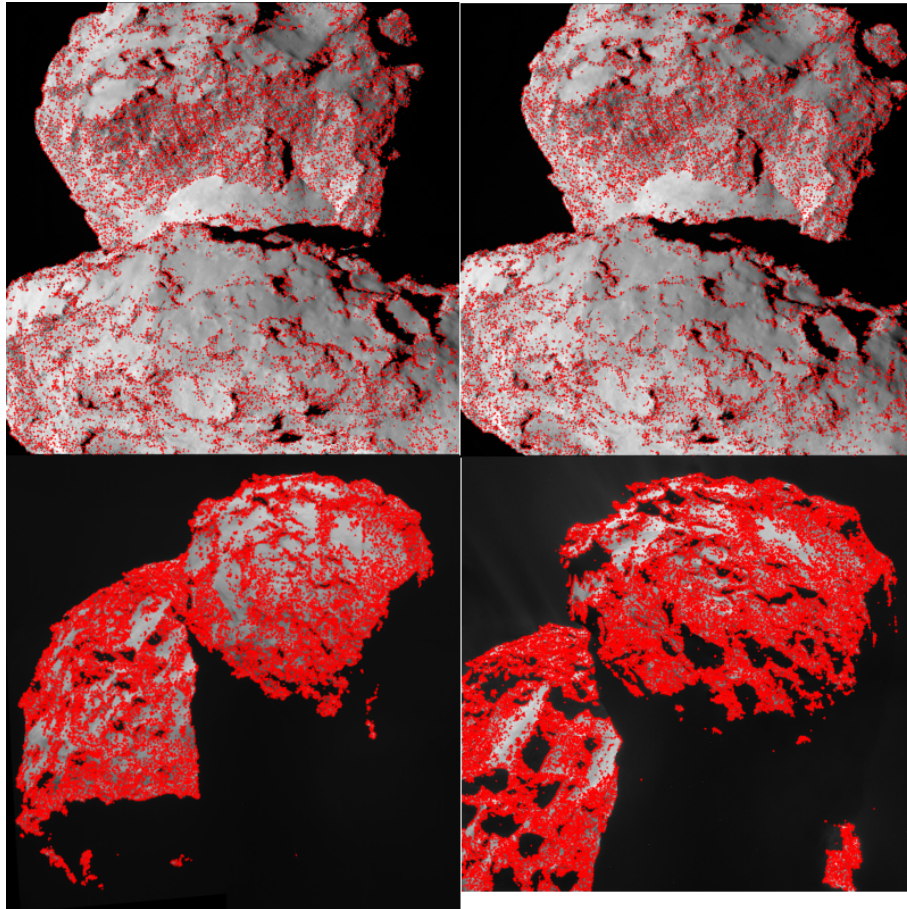


Figure. SIFT feature points

4.2 Robust Two-View Analysis

The objective of *Robust Two-View Analysis* is to find robust two-view features matches and accurate fundamental matrix estimation from SIFT features of two images.

We need to first obtain the *approximate nearest neighbor (ANN)*[5] matches of SIFT feature descriptors. For each set of descriptors, we construct a balanced kd-tree, which chooses split plane based on the variance on each dimension in the subtree. Then use *approximate* nearest neighbor search rather than *strict* for performance concern. Instead using a threshold of the Euclidean distance of two descriptor to indicate a match, we use a priority queue to match each feature, then use a threshold on the ratio of distance between the best match and the second. Features also need to be the ANN of *each other* to be accepted.

Then we use two-view epipolar constrain to stabilize the matches. Each pair of point correspondence must satisfy $x_1^T F x_2 = 0$, where F is the fundamental matrix. We exploit this attribute by combine RANSC and *normalized eight point algorithm* to filter outliers, and then estimate the *fundamental matrix* by SVD on the inliers then non-linear approximation by LMA.

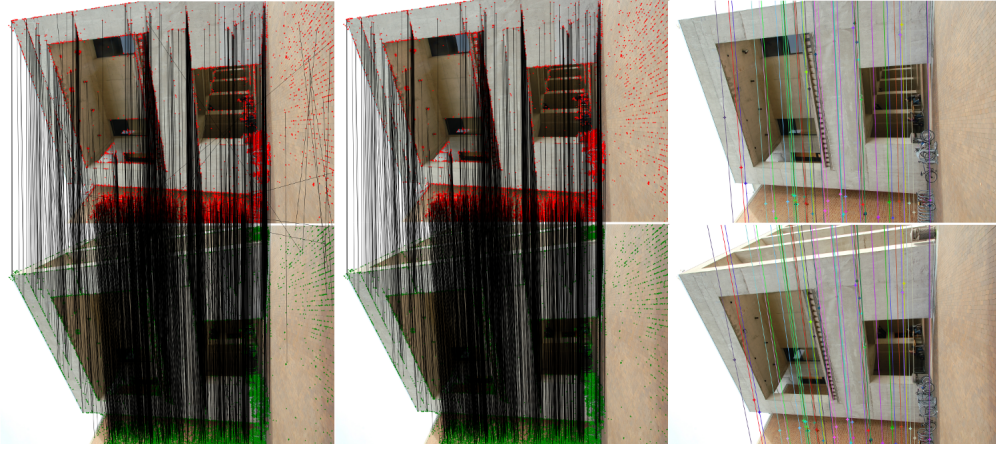


Figure 1. raw matches, robust matches and epipolar lines

4.3 Incremental Sparse Reconstruction

After the sparse structure is initiated, incremental recovery is performed. It will recover new cameras from recovered tracks, and triangulate more tracks, meanwhile use *SBA* and *Robust SBA* minimize projection error and keep the sparse structure robust, following is an overview of the algorithm:

Algorithm Incremental Sparse Reconstruction

```

tracks := input
recovered cameras := initial structure
recovered tracks := initial structure
candidates := un-recovered cameras view at least 20 recovered tracks
while candidates available
    selected := candidate view the most recovered tracks and those above 75 percentile
    new cameras := recover camera models of selected candidates using DLT
    recovered cameras += new cameras
    sparse bundle adjustment on new cameras and recovered tracks they observe
    new tracks := triangulate tracks at maximum view angle within recovered cameras
    recovered tracks += new tracks
    global robust sparse adjustment
    candidates := un-recovered cameras view at least 20 recovered tracks
end while

```

Recover new camera parameters is accomplished by obtain its projection matrix through DLT and then decompose it into (R, \mathbf{t}, K) , which can then be parameterized and refined.

When triangulate new points, big view angle is required. Only triangulate if the maximum view angle within recovered cameras is larger than a threshold.

4.4 Sparse Bundle Adjustment (SBA)

The objective of *Sparse Bundle Adjustment (SBA)*[3] is to minimize projection error of a sparse reconstruction by perform non-linear approximation on camera parameters and 3D points. It is a modified implementation of *Levenberg-Marquardt Algorithm* that take advantage of the sparse structure of multi-view geometry.

Sparse structure is exploited when calculating *Jacobian*, *damped Hessian* and δ_p . When solving

$$(H + d \cdot I) \cdot \delta_p = \epsilon$$

the damped Hessian can be split into $\begin{pmatrix} A & B \\ B^T & V \end{pmatrix}$, where V is a block diagonal matrix, and equation can be solved using Shur's component.

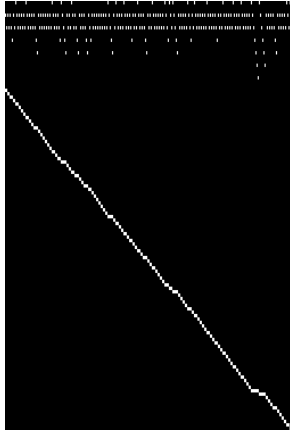


Figure 2. sparse jacobian matrix

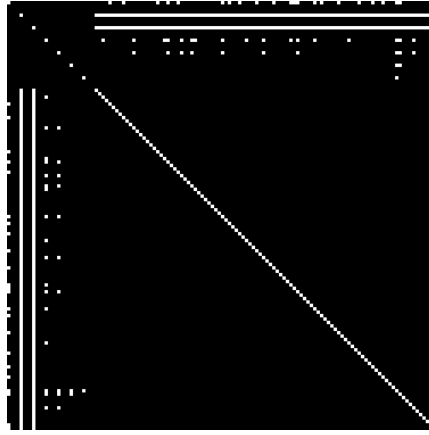


Figure 3. sparse hessian matrix

Robust SBA is SBA with outlier filter, it will refine the parameters and exclude outliers until no outliers can be found, following is an overview of the algorithm:

Algorithm Robust Sparse Bundle Adjustment

```

inliers := input tracks
cameras := input cameras
do
  sparse bundle adjustment on cameras and inliers
  for each cameravi in cameras:
    d80 := 80 percentile of reprojection error on cameravi
    thresholdvi := clamp(d80, 4, 16)
    outliersvi := tracks with reprojection error larger than thresholdvi
  end for
  outliers :=  $\bigcup$  outliersvi
  inliers := inliers - outliers
until outliers is empty

```


5 Application Framework

5.1 Parallelization

structure from motion pipeline contain large amount of independent calculations, which can be

We encapsulate computation intense logic into tasks, and put all tasks into a single script file **worker.js**, which can have more than one instance. Then **application.js** can instantiated it inside browser and uses WEBWORKER's event interface to call those tasks asynchronously. Task includes:

- SIFT task, take ImageData (as ArrayBuffer) and returns point buffer and vector buffer (as ArrayBuffer)
- Matcher task, take two camera shapes and two SIFT results, returns robust two view analysis
- Tracking task, takes SIFT point buffer of each image and all robust two-view analysis results, returns an array of tracks.
- Register task, takes tracks and fundamental matrices, returns a sparse point buffer and camera parameters

5.2 Data Model

Large amount of data flow through the pipeline, how to handle large data is especially critical inside browser.

A hierarchy of *ndarray*, *TypedArray*, and *ArrayBuffer* is used to handle large data, which uses *ArrayBuffer* as the underlying payload. For networking, it can be directly transfered through *XMLHttpRequest* as a native response type; for storage, it's supported in all implementations of *IndexedDB*; for multi-threading, it is a *TransferableObject* for *WebWorkers*. Data are stored in its minimum form through *TypedArray*, no JSON or text involved, which also means no parsing required.

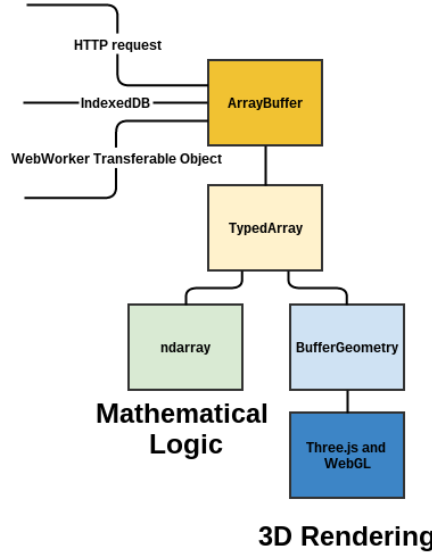


Figure 4. WebSFM data model

Following are the specifics of how data is represented as **ARRAYBUFFER**:

- SIFT features are stored in two separate parts, *points* and *vectors*. *Points* is a $4 \times l$ ndarray with an underlying **FLOAT32** typed array as buffer, which stores *row*, *column*, *orientation*, *scale* of each feature point. *Vectors* is a $128 \times l$ ndarray with an underlying **UINT8** typed array as buffer, which stores 128 dimensional descriptor of the cooresponding feature point in the *points* buffer. When matching features, only the *vectors* buffer is used, and when visualizing or tracking key points, only the *points* buffer is used.

- Surfels, both dense and sparse, is represented as two parts, *points* and *colors*. *Points* is a $3 \times l$ ndarray with an underlying FLOAT32 typed array as buffer, which stored the patch center coordinate. *Colors* is a $3 \times l$ ndarray with an underlying UINT8 typed array as buffer, which stores RGB color of the cooresponding patch in the *points* buffer.

Other data such as *camera paramters*, *feature matches* are insignificant in comparasion, so are directly stored as JSON objects.

To exam this approach, we take surfel model generated by PMVS which contains 450,000 patches. It costs 24.2MB in its original text form. this approach, data is stored in its minimum form, *points* and *colors* have combined size of 6.8MB.

5.3 Possible Improvement

WebGL is the door way to GPU acceleration inside browser, but haven't been fully exploited yet. The most time consuming tasks such as image convolution in SIFT and ANN matching can be GPU accelerated. This will be a future improvment.

5.4 Results

Bibliography

- [1] David G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, , 2004.
- [2] M. Bolitho M. Kazhdan and H. Hoppe. Poisson surface reconstruction. *Symp. Geom. Proc.*, , 2006.
- [3] Antonis A. Argyros Manolis I.A Lourakis. The design and implementation of a generic sparse bundle adjustment software package based on the levenberg-marquardt algorithm. , , 2004.
- [4] D Nistér. An efficient solution to the five-point relative pose problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, , 2004.
- [5] Nathan S. Netanyahu, Ruth Silverman, Angela Y. Wu Sunil Arya, David M. Mount. An optimal algorithm for approximate nearest neighbor searching in fixed dimensions. .
- [6] Jean Ponce Yasutaka Furukawa. Accurate, dense, and robust multi-view stereopsis. .