

3D Reconstruction Pipeline from Mobile Video (ORB-based)

1. Frame Extraction from Video

Frames are extracted from the handheld video stream (e.g. using OpenCV's `VideoCapture` or libraries like Decord). To avoid redundancy and speed up processing, one typically **downsamples** or selects keyframes. For example, taking every n th frame or using motion-based criteria ensures each pair has sufficient parallax. Keyframe selection can also use image similarity metrics or pose change thresholds ¹. On mobile, limiting resolution (e.g. resizing to 720p) and skipping near-duplicate frames is common. It's important to detect and **drop blurry or over/under-exposed frames** (e.g. via variance-of-Laplacian or histogram checks) to avoid bad feature matches.

Best practices: Convert frames to grayscale and optionally apply contrast enhancement (CLAHE or gamma correction) to mitigate lighting variation. For example, the AFE-ORB-SLAM pipeline uses adaptive gamma correction and unsharp masking to handle challenging illumination ². Pre-filtering (mild Gaussian blur) can reduce noise, but avoid heavy blur that removes corners.

Pitfalls: Too many frames (e.g. 60 FPS) yield large overlap with minimal motion; too few frames may miss fast motion. Large camera motion between frames can cause motion blur, so consider higher shutter speed or frame skipping. Uncalibrated rolling-shutter effects (common in phones) can distort fast motion – if present, either use short baseline or employ rolling-shutter models.

2. ORB Keypoint Detection & Description

ORB (Oriented FAST + Rotated BRIEF) is a fast, binary-feature extractor. In OpenCV Python, use `cv2.ORB_create(...)` with parameters like `nfeatures`, `scaleFactor`, and `nlevels` to control number of keypoints and scale pyramid. For mobile video, starting with ~1000–2000 features per frame is typical. ORB is *rotation-invariant* (it assigns an orientation to each keypoint) and partially scale-invariant via its image pyramid ³. This means moderate viewpoint changes are tolerated.

Parameter tuning: Increase `scaleFactor` or `nlevels` if expecting large scale changes. If ORB detects too few points in low-light, you can use an **adaptive FAST threshold**: AFE-ORB-SLAM, for example, boosts ORB performance in varying illumination by adjusting the FAST corner threshold and enhancing the image ². Using the Harris corner score (`ORB_create(scoreType=cv2.ORB_HARRIS_SCORE)`) can also weight stable corners more.

Robustness to motion blur/lighting/noise: ORB's binary descriptor (BRIEF) is sensitive to intensity noise but somewhat robust to monotonic brightness changes ³. In practice, normalize image brightness or use CLAHE to mitigate lighting shifts. Motion blur often reduces FAST keypoints; to counteract this, one can lower the FAST threshold or skip severely blurred frames. Direct methods (like LSD-SLAM) are noted to be

more robust to blur than feature-based ones ⁴, so blending approaches (see Section 7) can help in blur. Noise can be reduced by denoising filters; however, too much filtering will remove corner features.

Code example:

```
orb = cv2.ORB_create(nfeatures=1500, scaleFactor=1.2, nlevels=8)
kp, des = orb.detectAndCompute(gray_frame, None)
```

(See e.g. a tutorial on ORB usage ⁵ ⁶.) This yields keypoints and binary descriptors.

3. Feature Matching Between Frames

For each pair of frames (usually consecutive or keyframe-to-keyframe), match ORB descriptors. Using OpenCV in Python: initialize a brute-force matcher with Hamming distance (`cv2.BFMatcher(cv2.NORM_HAMMING)`) and perform k-NN matching ($k=2$). Apply **Lowe's ratio test** (e.g. retain match if best distance $< 0.75 \times$ second-best) to reject ambiguous matches ⁷. It's also advisable to enforce **cross-check** (match $A \rightarrow B$ and $B \rightarrow A$) or an additional geometric filter: for instance, the LearnOpenCV pipeline also checks that the matched keypoints are not too far apart in normalized coordinates ⁷.

After initial matching, use RANSAC to remove geometric outliers. For example, compute the fundamental or essential matrix with RANSAC and keep only the inliers. The `match_frames` routine from a Python SLAM example does this: it matches ORB features with ratio test, filters by a small Euclidean distance, then RANSAC-fits a Fundamental matrix to retain inliers and obtains the relative pose ⁷ ⁸. In practice, you can just pass matches to `cv2.findEssentialMat()` with RANSAC.

Optimizations: If matching is slow, consider FLANN with LSH for binary descriptors (OpenCV's `FlannBasedMatcher` with `index_params={'algorithm':6}`), or track keypoints via optical flow (LK) between neighboring frames to reduce descriptor matching. Also, you can spatially restrict matching (e.g. epipolar line search after an initial pose estimate) to speed up dense matching.

Pitfalls: Repetitive textures (walls, sky) can cause many false matches. Use a low ratio threshold or a higher `uniquenessRatio` if using SGBM. Very small motions yield few disparities and may confuse matching; ensure enough baseline (skip nearly identical frames). Incorrect matches will later degrade pose estimation, so always use RANSAC thoroughly.

4. Camera Pose Estimation (Essential Matrix)

With matched inlier points (after RANSAC), estimate the relative camera pose between the two frames. First ensure you have or assume camera intrinsics (focal length f and principal point). For phones, f can be obtained via camera specs or estimated (common approach: use image width, e.g. $f \approx 0.8 * \text{width}$ if unknown).

Use `cv2.findEssentialMat(pts1, pts2, cameraMatrix, method=cv2.RANSAC, prob=0.999, threshold=1.0)` to compute the essential matrix **E**. A robust threshold (in pixels) is important: too low

and true inliers get dropped; too high and many outliers slip through. After obtaining \mathbf{E} , recover the relative rotation \mathbf{R} and translation \mathbf{t} (up to scale) via `cv2.recoverPose(E, pts1, pts2, cameraMatrix)`. This yields \mathbf{R} and a unit \mathbf{t} vector. Internally, `recoverPose` chooses the correct among the four possible solutions by enforcing positive depth.

Accuracy and limitations: Pure rotation yields an indeterminate \mathbf{t} (no parallax), so the essential matrix becomes rank-deficient. Similarly, if points lie on a plane, pose recovery is ambiguous. RANSAC helps avoid degenerate cases, but do not estimate pose if inlier count is too low (e.g. < 50) or if motion is nearly rotational. The recovered translation has an unknown scale; for metric depth you need known baseline or an external scale (e.g. known object size). For multiple frames, one usually accumulates the scale by triangulating points or using an external scale reference.

Accuracy considerations: The pose estimate can be noisy due to image noise or mismatches. One improvement is to refine the pose by **bundle adjustment (BA)** over multiple frames and 3D points (e.g. `g2o` or `Ceres`), as done in ORB-SLAM. In Python, small-scale BA can be done with libraries like `g2o` or `ceres-python`. At a minimum, one can reproject matched points and compute reprojection error to gauge quality.

Pitfalls: Fast motion blur or rolling shutter can violate the pinhole model. If the phone provides inertial data (IMU), tightly-coupled vision-inertial methods can improve pose. Without external IMU, be cautious: false matches can flip the pose. Validate that inliers produce a coherent motion.

5. Stereo Rectification of Frame Pairs

To prepare for dense matching, rectify each chosen frame pair so corresponding epipolar lines align horizontally. Two cases:

- **Calibrated (intrinsics known):** Given \mathbf{R} , \mathbf{t} , and camera matrices, use OpenCV's `cv2.stereoRectify(K, dist, K, dist, imageSize, R, t, R1, R2, P1, P2)` to compute rectification transforms. Then use `initUndistortRectifyMap` + `remap` to warp images to a common plane. This ensures disparities are horizontal shifts.
- **Uncalibrated (baseline/ints unknown):** Compute the fundamental matrix \mathbf{F} (via RANSAC if needed). Then use `cv2.stereoRectifyUncalibrated(points1, points2, F, imageSize)` to get homographies \mathbf{H}_1 , \mathbf{H}_2 for rectification⁹. This does not require \mathbf{K} , but assumes your point correspondences cover the scene well. After getting $\mathbf{H}_1/\mathbf{H}_2$, warp images by `cv2.warpPerspective`.

The latter method (uncalibrated rectification) is useful if you lack exact intrinsics. The OpenCV docs note that it “does not need to know the intrinsic parameters”¹⁰, but it “heavily depends on the epipolar geometry.” It’s wise to undistort images first if lens distortion is significant¹¹.

Practical tips: Choose frame pairs with enough baseline (e.g. skip several frames) to ensure a discernible disparity. Very small baseline yields almost no disparity (flat plane). After rectification, verify alignment by checking a few corresponding points lie on the same rows.

Pitfalls: If the baseline is extremely small or the motion is mainly forward (depth far away), rectification yields very low disparity (depth uncertainty). If the scene has strong depth discontinuities, block-based matching can struggle (see next section). Also, uncalibrated rectification assumes planar homography – for wide-FOV cameras with distortion, residual misalignment may occur unless undistortion is applied ¹¹.

6. Dense or Semi-Dense Depth Estimation

Once images are rectified, obtain depth (disparity) via either classical stereo or learning methods:

- **Classical stereo (Block Matching / SGBM):** Use OpenCV's `StereoBM` or `StereoSGBM` to compute a disparity map from the rectified pair. StereoBM is faster but less accurate; StereoSGBM (semi-global block matching) yields smoother, more accurate maps at higher cost. Tuning parameters (numDisparities, blockSize, uniquenessRatio, speckleWindow, etc.) is crucial – typically start with blockSize=5–9, numDisparities=multiple of 16 depending on scene depth range ¹². The result is a *disparity map* that can be converted to depth via $Z = fB/d$ ¹³ (where d is disparity in pixels, B is baseline, f is focal).
- **Triangulation:** With projection matrices from `stereoRectify` (P_1, P_2), one can also use `cv2.triangulatePoints(P1, P2, pts1, pts2)` on matched features to get sparse 3D points. This supplements the disparity map with metric scale if B and f are known (baseline can be recovered up to scale from extrinsics).
- **Multi-view stereo (MVS):** If multiple overlapping views are available (more than two), methods like *plane-sweep stereo* or *multi-frame SGM* improve robustness. For example, the Mobile3DRecon system employs multi-view SGM: they aggregate cost volumes across several frames via a weighted census transform and perform winner-take-all to get an initial depth ¹⁴, then refine it with a neural network.
- **Deep learning (monocular):** Single-image depth networks (e.g. **MiDaS**, **DPT**) can predict relative depth from one photo. MiDaS v3 provides models ranging from heavyweight (BEiT, Swin transformers) to lightweight (Swin2-Tiny, LeViT) ¹⁵. Notably, MiDaS v2.1 introduced a **mobile-optimized model** that runs in real-time on phone CPUs ¹⁶. These methods output depth up to an unknown scale; one can align scale by known objects or additional stereo if available.
- **Deep learning (stereo):** Learned stereo networks (e.g. PSMNet ¹⁷, StereoNet, RAFT-Stereo) take rectified image pairs and output dense disparity. They usually require GPU acceleration and training on synthetic/real datasets. PSMNet, for instance, uses pyramid pooling and 3D convolutions to exploit context ¹⁷. Modern systems like **HSM** or **PSMNet** achieve high accuracy but may not run on mobile in real time. For lighter inference, some networks use 1D convolutions or quantization.
- **Sparse-to-semi-dense:** Classical sparse ORB yields few points. To densify, one can propagate depths from reliable pixels (e.g. along edges) or fuse multiple views. Direct SLAM systems like LSD-SLAM recover a semi-dense depth map on image edges by photometric tracking ⁴. In our pipeline, one strategy is to compute depth only for pixels where stereo confidence is high (e.g. edges), yielding a *semi-dense* map. Alternatively, reproject disparities to a point cloud and apply surface reconstruction or volumetric integration for a mesh.

Depth Filtering: Stereo disparity maps often contain outliers in textureless or specular regions. It's common to apply median or bilateral filtering, or occlusion inpainting. The OpenCV function `reprojectImageTo3D` can convert disparity to 3D points (mask invalid disparities). For semi-dense output, one can threshold by disparity gradient or matching cost. Confidence estimation (e.g. based on left-right consistency) helps mask unreliable areas.

Pitfalls: Stereo matching fails on texture-less areas, occlusions, or large radiometric differences (lighting changes between views). Active filtering of such areas (leaving holes) may be preferable to wrong depths. Monocular networks suffer from scale ambiguity and may produce less detailed geometry. Semi-dense methods like LSD-SLAM avoid matching in flat regions by design ⁴ but require consistent lighting.

Below is a summary table of common depth methods:

Method	Input	Output	Notes
Block Matching (StereoBM) ¹⁸	2 rectified images	Disparity map	Fast, simple; good for real-time CPU. Struggles with low texture/noise.
Semi-Global Matching (SGBM)	2 rectified images	Disparity map	Higher quality than BM; more tuning (block size, speckle filters).
PSMNet (CNN) ¹⁷	2 images	Disparity map	CNN-based (stacked hourglass) stereo; high accuracy, needs GPU.
MiDaS (monocular) ¹⁶	1 image	Relative depth	State-of-art monocular depth; lightweight models available for mobile.
DPT (monocular)	1 image	Relative depth	Vision-transformer depth; high detail.
LSD-SLAM (direct) ⁴	Monocular video	Semi-dense depth	Semi-dense depth on high-gradient regions; robust to blur/low-texture ⁴ .

7. From Sparse ORB Points to Semi-Dense 2.5D Map

ORB-based matching and triangulation yield a **sparse point cloud**. To achieve a richer “2.5D” reconstruction (a depth map or semi-dense surface), one can fuse information:

- **Multi-view integration:** Accumulate depth estimates from multiple frame pairs (via SGM or networks) into a single model. For example, the Mobile3DRecon pipeline incrementally integrates each keyframe's depth map (after filtering) into a global mesh using TSDF fusion ¹⁹ ²⁰. On mobile, one might select a few well-spaced keyframes, compute their depth maps (using stereo or monocular methods), and merge them.
- **Edge-based densification:** Use detected ORB points to seed depth along edges. One approach is to run small-window photometric stereo (like in LSD-SLAM): track image gradients over several frames to solve depth on edges. This yields a semi-dense map without explicit stereo. The LSD-SLAM system

uses exactly this idea ⁴, though implementing it requires direct image optimization rather than feature matching.

- **Interpolation:** After obtaining a dense disparity map from stereo (as above), use it directly as a 2.5D depth map. If only sparse points are available, one can use interpolation techniques (e.g. guided image filtering of the sparse depth) to create a semi-dense map.

Robustness under uncontrolled lighting/noise: Illumination changes can disrupt photometric methods. If using stereo, ensure the two images are captured with similar exposures. Using gradient-based costs (Census transform in SGM) helps resist monotonic brightness changes. For monocular networks, diversity in training data (MiDaS mixes datasets) gives some robustness. Always apply post-filtering (median, bilateral) to smooth out noise spikes.

Pitfalls: Dynamic objects in the scene will produce inconsistent depth across frames; filter them by temporal consistency or mask moving regions. Large brightness changes (sun glare) can create artifacts in stereo; focus on textured regions. Finally, scale drift remains a challenge: if monocular depth is used, scale can drift without a reference.

8. Code and Implementation Resources

Several open-source projects and tutorials illustrate parts of this pipeline in Python:

- **SLAMPy:** A Python-based monocular ORB-SLAM example that uses OpenCV ORB, Lowe's ratio test, and g2o for bundle adjustment ²¹. It demonstrates keypoint detection, matching, pose estimation, and mapping (using Pangolin for visualization).
- **PyOrbSlam:** A Python wrapper/examples around ORB-SLAM (using OpenCV and the g2o optimizer) on GitHub ²².
- **LearnOpenCV Monocular SLAM tutorial:** A step-by-step guide in Python covering ORB, BFMatcher, essential matrix, and bundle adjustment (see the "match_frames" function using ratio test and RANSAC ⁷ ⁸).
- **OpenCV sample code:** The OpenCV documentation and forums have examples of using ORB (`cv2.ORB_create`), `BFMatcher`, `cv2.findEssentialMat`, `recoverPose`, `cv2.stereoRectifyUncalibrated`, and `cv2.StereoBM/StereoSGBM` ⁹ ¹⁸.
- **Depth estimation repos:** The MiDaS GitHub provides code/models for monocular depth (PyTorch) ¹⁶. The PSMNet repository (Chang & Chen, CVPR 2018) offers a PyTorch stereo network implementation ²³.

9. Common Pitfalls and Optimizations

- **Frame issues:** Motion blur or very low light can cause few ORB features. Consider denoising or skipping such frames. Bright flicker (from auto-exposure) can break descriptor matching; fixing exposure or using exposure metadata to normalize images can help.
- **Repetitive textures:** Uniform areas yield few matches, so rely more on edges. You might combine ORB with other detectors (e.g. Harris corners) if needed.
- **RANSAC tuning:** If too many inliers are rejected, increase the inlier threshold (especially on high-res images). Conversely, too many outliers will corrupt pose. Visualize inlier matches to verify.
- **Scale ambiguity:** Monocular pipelines (Essentials) have unknown scale. Without external metric info (like known object size or IMU), the reconstruction is up to a scale factor. To recover scale, one can

calibrate movement (e.g. with AR markers) or assume phone's focal length (getting metric depth from disparities).

- **Performance:** For real-time, reduce feature count or use GPU-accelerated algorithms (OpenCV CUDA stereo, or mobile NN with TensorRT). On phones, smaller models (MiDaS-Tiny ¹⁶) or optimized C++ code can boost speed.
- **Alternative approaches:** Instead of ORB, one could use other features (AKAZE, BRISK) or even direct methods if lighting is stable. For dense mapping, SLAM systems like ORB-SLAM3 (with RGB-D or stereo modes) provide robust pipelines but are C++-based.

10. Conclusions

Building a 3D reconstruction pipeline from handheld video involves careful feature engineering (ORB tuning), robust matching (ratio tests, RANSAC), and multi-step estimation (pose via essential matrix, rectification, depth). ORB features offer a good speed/accuracy trade-off on mobile, but require safeguards (adaptive thresholds, pre-filtering) under challenging conditions ² ⁴. Depth can be obtained via traditional stereo or modern CNNs; using them in tandem (classical SGM + CNN refinement ¹⁴ ²⁴) is a strong strategy for mobile devices. Finally, converting sparse points to a richer 2.5D model may use multi-view fusion or direct edge-based methods. By combining these steps with the cited best practices and references, one can assemble an end-to-end pipeline suited for smartphone video input.

Sources: We leveraged open-source documentation, recent research (e.g. ORB-SLAM, MiDaS, Mobile3DRecon) and community code examples to compile these guidelines ³ ⁷ ⁹ ¹⁴ ² ¹⁶ ¹⁷ ⁴. These cover both theoretical foundations and practical tips for each stage of the pipeline.

- 1 Keyframe Selection for Visual Localization and Mapping Tasks: A Systematic Literature Review
<https://www.mdpi.com/2218-6581/12/3/88>
- 2 (PDF) AFE-ORB-SLAM: Robust Monocular VSLAM Based on Adaptive FAST Threshold and Image Enhancement for Complex Lighting Environments
https://www.researchgate.net/publication/360656953_AFE-ORB-SLAM_Robust_Monocular_VSLAM_Based_on_Adaptive_FAST_Threshold_and_Image_Enhancement_for_Complex_Lighting_Environments
- 3 4 ORB-SLAM: A Versatile and Accurate Monocular SLAM System
https://courses.cs.washington.edu/courses/csep576/21au/resources/ORB-SLAM_A_Versatile_and_Accurate_Monocular_SLAM_System.pdf
- 5 6 7 8 Understanding Monocular SLAM implementation in Python OpenCV
<https://learnopencv.com/monocular-slam-in-python/>
- 9 10 11 OpenCV: Camera Calibration and 3D Reconstruction
https://docs.opencv.org/4.x/d9/d0c/group__calib3d.html
- 12 18 Stereo Camera Depth Estimation With OpenCV (Python/C++)
<https://learnopencv.com/depth-perception-using-stereo-camera-python-c/>
- 13 17 23 Pyramid Stereo Matching Network
https://openaccess.thecvf.com/content_cvpr_2018/papers/Chang_Pyramid_Stereo_Matching_CVPR_2018_paper.pdf
- 14 19 20 24 Mobile3DRecon: Real-time Monocular 3D Reconstruction on a Mobile Phone
<https://zju3dv.github.io/mobile3drecon/>
- 15 16 GitHub - isl-org/MiDaS: Code for robust monocular depth estimation described in "Ranftl et. al., Towards Robust Monocular Depth Estimation: Mixing Datasets for Zero-shot Cross-dataset Transfer, TPAMI 2022"
<https://github.com/isl-org/MiDaS>
- 21 GitHub - Akbonline/SLAMPy-Monocular-SLAM-implementation-in-Python: Pythonic implementation of an ORB feature matching based Monocular-vision SLAM.
<https://github.com/Akbonline/SLAMPy-Monocular-SLAM-implementation-in-Python>
- 22 GitHub - mathbloodprince/PyOrbSlam: Orb feature based SLAM (simultaneous localization and mapping) using OpenCV for feature extraction and pose estimation, g2o for bundle adjustment, and OpenGL for visualization.
<https://github.com/mathbloodprince/PyOrbSlam>