

Notes on DSO

Xiang

September 28, 2018

1 Introduction on DSO

1.1 What is DSO

DSO (Direct Sparse Odometry) is a visual odometry system published in 2016 by Dr. Jakob Engel from the Computer Vision Laboratory at the Technical University of Munich (TUM) [2]. For the laboratory home page, see <https://vision.in.tum.de/research/vslam/dso>. The code link is here: <https://github.com/JakobEngel/dso>. For now there are many different versions of DSO like stereo DSO, IMU-intergrated DSO, DSO with loop closing, DSO with auto photometric calibration, DSO with rolling shutter camera, etc. Some of them are from our CV group and the others are from public community. Unfortunately our group only releases the basic DSO version, and if you are interested on stereo DSO or IMU-DSO, please look for other open-source versions in the website, like https://github.com/HorizonAD/stereo_dso.

DSO is already well-explained in the paper of [2], however the code of DSO is not very clearly-written compared to other open-source systems like ORB-SLAM2 and SVO (thanks to Jakob), which makes it very difficult for researchers (and me) to conduct follow-up research based on it. This note is to explain some of the basic knowledges about DSO's code, like how it is organized, what is the working pipeline of DSO, etc.

In order to understand this note you should at least have the basic knowledge of visual SLAM, like 3D rigid-body motion and bundle adjustment. You'd better read the DSO's paper at first before reading this material.

1.2 Direct method

DSO is a keyframe visual odometry based on direct method. It is **not** a complete SLAM system and does not include loop closing detection or map reusing, which means it *won't* reuse the previously built map. DSO keeps a small window of active keyframes (sometimes called **Sliding Window Filter** (SWF) algorithm), Once a keyframe is out of the active window, you cannot use its information anymore (which is called as **Marginalization** and we will explain it later). Therefore, it inevitably has accumulative errors (on rotation, translation, and scale), which are small but cannot be eliminated.

DSO is one of the few visual odometry systems using fully direct method. In contrast, SVO [4] is a semi-direct method which use direct tracking as the front-end but still minimizes the geometry error in back-end, and LSD-SLAM [3] is only a frame-by-frame odometry with out a back-end¹. DSO also makes the direct tracking much more robust than SVO and LSD-SLAM, by introducing the photometric parameters (exposure time, affine light transform) and a carefully tweaked optimization process. The direct method has a very important advantage (but usually ignored) compared to the traditional feature point method: it puts data association and pose estimation into an unified nonlinear optimization problem, while in the feature method, they are solved step by step, that is, the data association is firstly obtained by matching the feature points, which is a discrete problem and cannot be solved together with pose estimation. Then it estimates the pose based on the given and fixed association. These two steps are usually independent. In the second step, the re-projection error can be used to judge the outliers in matching, but we normally don't do re-matching based on the pose estimation (theoriotically, we still can do it like [1] but it is expensive in computation).

Feature matching is actually a very instable step. First you need to have enough corners (points on edges and areas cannot be accurately matched). Second, these corners should be distinctive when comparing with each other.

¹You can simply consider DSO as LSD plus a photometric BA.

For example, you should not extract corners in a brick wall or stone road because they all look the same. This makes the SLAM system not robust in places where you don't have such distinctive corners, and unfortunately this is usually the case of real-world environments. But in direct method, you don't need those matched points. You can make use of corners, edges, circles, stars and even smooth areas like the wall, because we don't need to previously determine the matching result before estimating the camera pose. This unified optimization framework in direct method, although more complicated than feature matching, is much more robust in low-texture areas.

So in DSO, we don't have the concept of "matched points." Each 3D point is derived from a host frame, multiplied by the depth value and then projected to another target frame, creating a reprojection residual. As long as the residuals are within some reasonable limits, they are considered to be projected by the same point. From the perspective of data association, there is no such thing like fixed $a_1 - b_1$, $a_2 - b_2$ in this process, and there may be cases like $a_1 - b_1, a_2 - b_1, a_3 - b_1$, but DSO does not care. Again, as long as the residuals are not large, we will just treat it as the same point. This is a very important part in direct method. In the feature point method, we can record in which frames a map point is seen, and even find the image descriptors in each frame; but in DSO, we will try to project each point into **all** active frames (unless they are out-of-boundary or occluded), calculate its residual in each frame, and do not care about the one-one correspondence between points.

In the back end, DSO builds a sliding window consisting of several (5-7) keyframes. This window persists throughout the VO process and has a set of methods to manage the addition of new data and the removal of old data. The front-end tracking part will determine whether the new frame should be inserted into the back-end as a new key frame by certain conditions. At the same time, if the number of key frames found by the backend is larger than the window size, one of the frames will be removed according to some specific conditions. Please note that the removed frame is not necessarily the oldest frame on the timeline. There are some complicated conditions to decide which one should be marginalized out. And once a key-frame is removed, the information within this key-frame will be **marginalized** into the active window, and become the **prior** of the window.

Since the direct method requires comparison of image information, we need to store at least the images of the keyframes, which makes the window size cannot be too large. Besides, image functions are not smooth, and are easily disturbed by lighting condition. Therefore, DSO proposed the photometric calibration, which can make the direct method more robust after calibrating the camera's exposure time, vignetting angle and gamma response. For cameras that are not photometrically calibrated, the DSO also dynamically estimates the photometric parameters (specifically, an affine change parameter, denoted as a and b). And there are also automatic photometric calibration for DSO described in a not-published paper in our group. This process models the imaging process of the camera, so it will perform better for image changes caused by different camera exposures. However, this cannot deal with the changes from the light source or reflections.

OK now let's talk about the (somewhat terrible) DSO's code. For more details about the theory, please take a look at the original DSO's paper.

2 About the Code

2.1 Framework of DSO

First let's look at the general framework of DSO. I removed some unimportant classes and structures for the readers. See Fig.2.1 for details (English version will be added later).

The DSO code consists of four parts: the system and each algorithm are integrated in src/FullSystem, and the back-end optimization is located in src/OptimizationBackend, which constitutes most of the core content of DSO. Src/Utils and src/IOWrapper are some undistorting, dataset reading and writing and visual UI code. Let's look at the core part of FullSystem and OptimizationBackend.

As shown in the statements above, in FullSystem, DSO will maintain a sequence of keyframes inside a sliding window. The data of each frame is stored in the FrameHessian structure, and FrameHessian is a frame with state variables and Hessian information. Then, each frame also carries some map point information, including:

1. point Hessians is information about all active points. The so-called active points mean that they are in the field of view of the camera, and the residual term is still participating in the calculation of the optimization part;
2. point HessiansMarginalized is a map point that has been marginalized;
3. point HessiansOut is a map point that is marked to be an outlier;

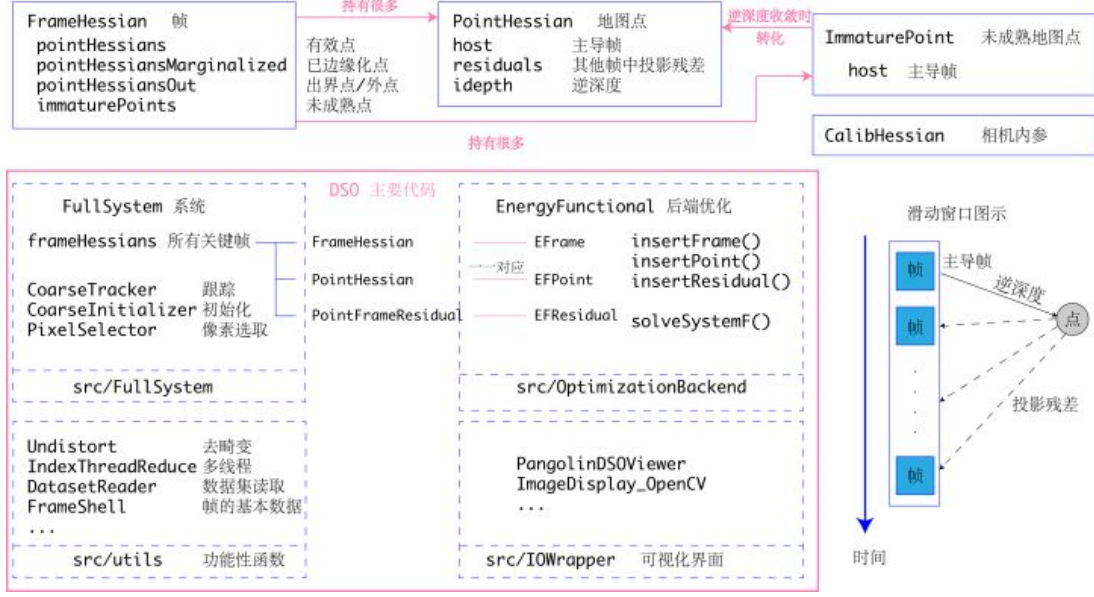


Figure 1: Framework of DSO.

4. And immaturePoints is the information of immature map points.

In monocular SLAM, all map points have only one 2D pixel coordinate at the beginning. The depth is unknown. This point is called an immature map point in the DSO and stored in the ImmaturePoint class. As the camera moves, DSO tracks these immature map points on each image. This process is called **Trace** and is actually a process of searching along the epipolar line, very similar with svo's depth filter. The **Trace** process determines the inverse depth of each Immature Point and its range of variation. If the depth of the Immature Point (actually the inverse of the depth, ie the inverse depth) converges in the process, then we can determine the three-dimensional coordinates of this immature map point, forming a normal map point. A map point with three-dimensional coordinates, called PointHessian in DSO. In contrast to FrameHessian, PointHessian also records the three-dimensional coordinates of this point, as well as Hessian information.

Unlike many other SLAM schemes, DSO uses a single parameter to describe a map point, its inverse depth. For most programs, such as ORB-SLAM, the coordinates of x , y , and z of the map point are recorded. The inverse depth parameterized form has the advantages of simple form, similar Gaussian distribution, and more robustness to distant scenes [5], but based on Bundle adjustment of inverse depth parameterization, each residual term needs to calculate one more than the usual BA. Jacobian matrix. In order to use the inverse depth, each PointHessian must have a host frame indicating that the point is back-projected from the frame.

Thus, all information of the sliding window can be described by several FrameHessian, plus the PointHessian carried by each frame. All PointHessian can be projected in any frame from a host frame to create a residual term, which is recorded in PointHessian::residuals. Adding all the residuals constitutes an optimization problem that the DSO needs to solve. Of course, due to motion and occlusion, not every point can be successfully projected to all of the remaining frames, so we also need to set the state of each point as: valid /marginal/invalid. Points of different states are stored in the three containers of pointHessians/pointHessianMarginalized/PointHessiansOut of its host frame.

In addition, DSO takes into account the camera's internal parameters, exposure parameters and other information as optimization variables. The camera internal parameters are expressed by the pinhole camera parameters f_x , f_y , c_x , c_y , and the exposure parameters are described by two parameters a , b . Please read the photometric calibration paper for details.

The backend optimization part has separate Frame, Point, and Residual structures. Since the optimization goal of DSO is to minimize energy (similar to cost in Ceres), the classes related to the backend start with EF and correspond one-to-one with the instances stored in FullSystem, holding each other's pointers. The optimization part is managed by the EnergyFunctional class. It takes all the frames and points from FullSystem, optimizes them, and returns the optimization results. It also contains all the frame and point information in the entire sliding window and is responsible

for handling the actual nonlinear optimization matrix operations.

2.2 Changes of LDSO

LDSO has made a lot of changes to the framework of DSO, in order to make the code structure more clear. The code structure of LDSO is mainly composed of the following folders:

1. include and src: header files and source code files.
2. examples: some examples, such as running on Kitty or on EUROC.
3. thirdparty: third-party libraries, sophus, dbow3 and g2o.
4. vocab: a dictionary for loop detection

The source code is divided into several subfolders:

1. In include/: the basic code structure, like Frame, Point, Feature, etc.
2. include/frontend is the front-end related code, including the feature detection, tracking, as well as the system interface.
3. Under include/internal is the backend and some internal algorithms, basically the content of the DSO, but for some reason I do not want to modify them.
4. Under include/internal/OptimizationBackend is the backend that DSO implements itself. Since the general optimization libraries like Ceres or g2o do not natively support incremental optimization, DSO implements one by itself. This backend uses a lot of hacks², making it hard to read and not general, but faster than the others.

In the LDSO I extracted the main structure and made the whole logic look clearer. In general, in LDSO you will have some **Frames**, and these **Frames** will see some **Features**. If the depth of this **Feature** is correctly estimated, it will produce a corresponding **Point**. This structure is similar to SVO and is simple enough. So, if you want to get the entire track trajectory or save the entire map, just iterate through all the **Frames** and **Points** and save them.

These classes are the simplest data structures of the program, of course, we still need to store some extra data in the actual operation. For example, for any **Frame**, I want to estimate its pose, then there may be an estimated state, a previously estimated state (for backup), and an increment. I put these in the internal/**FrameHessians**, similar for **Point** to internal/**PointHessian**. In DSO, the old **Frame** and **Point** are immediately destructed from the inside, but in LDSO, I need to save them for rebuilding the map after loopback. So you can access the state of each **Frame** and **Point** at any time. **Frame** will hold a **FrameHessian** pointer and deconstruct it if not needed.

LDSO also made some simplifications on the backend. In DSO, the backend is implemented by the EnergyFunctionals class. It uses some of its own EFPoints, EFFrame which have some overlap with the front end, so you often see a state data is passed from the front end to the backend, and passed back when calculation is finished. LDSO directly lets **EnergyFunctional** handle **FrameHessian** and **PointHessian**, and the results are stored directly in these classes, making the backend structure simpler.

At the algorithm level, LDSO has a feature extraction and matching process, in the frontend/FeatureDetector, frontend/FeatureMatcher, to achieve the process of matching and extraction. Frontend/LoopClosing implements loopback detection. If the loopback is detected, then one of the Map classes completes the optimization of Sim3.

LDSO saves the front/back position of the loopback for each frame. The optimized pose should be uniform in scale. Finally, the map corresponding to the optimized pose is displayed in the visualization window. I admit that the backend of DSO is difficult to change. I think the opinions of others in our group are the same.

2.3 Tracking

The frontend pipeline of DSO is illustrated in Fig.2.3. Although the specific meaning of each step is not explained, the general flow of the DSO can be seen. From the above figure, we can briefly summarize the behavior of DSO:

²SSE acceleration of computing accumulated Jacobians, sharing some internal computation results, etc.

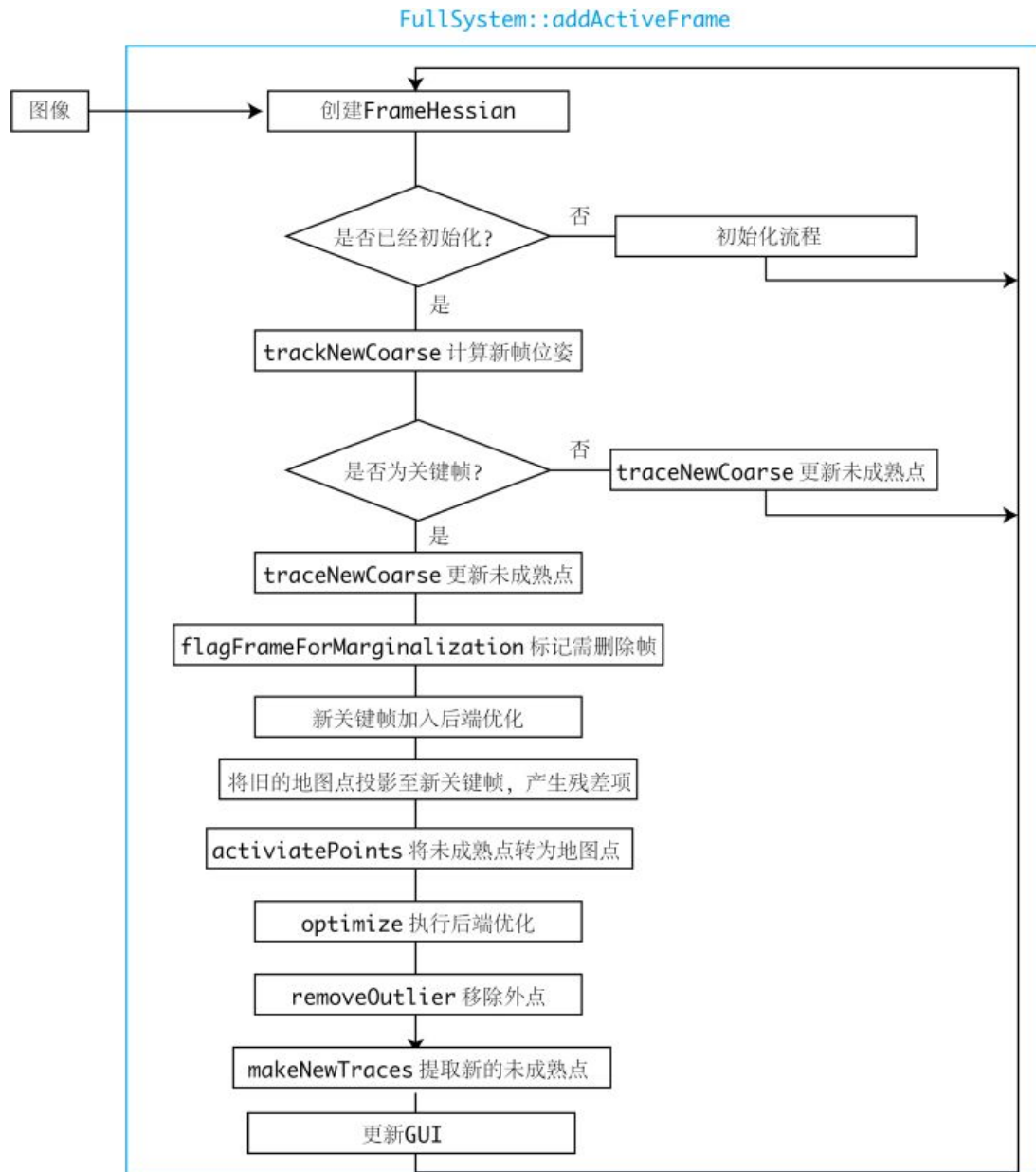


Figure 2: Pipeline of the frontend.

1. For non-keyframes, the DSO only calculates its relative pose with the last keyframe, and updates the depth estimate for each immature point with the frame image;
2. The backend only tackles the optimization of keyframes. In addition to some memory maintenance operations, the main processing for each key frame is: adding new residual items, removing the wrong residual items, and extracting new immature points.
3. The entire process is in one thread, but there may be multi-threaded operations inside.

2.4 Backend

2.4.1 Marginalization

The keyframes and the residuals of the map points, form the content of the entire sliding window. To optimize these frames and points, we will iterate using the Gauss-Newton or Levenberg-Marquardt methods. In the iteration, all residual terms can be combined into a large linear equation:

$$\mathbf{J}^T \mathbf{W} \mathbf{J} \delta \mathbf{x} = -\mathbf{J}^T \mathbf{W} \mathbf{r}, \quad (1)$$

This process is also called Marginalization, when we marginalize all the points and put their information into the pose part, forming a priori of the camera pose. This part of the knowledge should be familiar to every SLAM researcher.

So where is the use of marginalization in DSO?

First of all, DSO's BA, like the traditional BA, has the above steps. Therefore, when the DSO solves the BA, it marginalizes the information of all the points and calculates the optimized update amount. However, unlike the traditional BA, the upper left part of the DSO, ie \mathbf{B} in the formula, is not a diagonal block but a priori. In the traditional BA, this part is a diagonal block. The main reason is that the prior motion of the camera movement is not known, and the sliding window of the DSO is calculated by a certain means. The a priori here comes mainly from two parts: When a point is edged, a priori is generated between the common frames of the point; When a frame is edged, a priori is generated between other frames in the window; Here, "marginalization", the specific operation is the same as the marginalization mentioned above. That is to say, through Shure complement, use part of the matrix to eliminate another part of the matrix. However, the meaning of the actual operation is different. In the marginalization of BA, we hope to use marginalization to accelerate the solution of the whole problem, but after solving the problem, these frames and points still exist in the window! The marginalization in the sliding window means that we no longer need this point/this frame. When it is marginalized, we pass its information to the subsequent a priori, and will not use this point/this frame anymore! The reader must be clear about this difference, otherwise you will encounter problems during the understanding process. We might as well call the latter "permanent marginalization" to distinguish. So how does DSO permanently marginalize a frame or point? It follows the following guidelines: If a point is no longer in the field of view of the camera, marginalize this point; If the number of frames in the sliding window has exceeded the set threshold, then select one of the frames for marginalization; When a frame is marginalized, the map points that dominate it will be removed and will not participate in future calculations. Otherwise this point will form a priori with other points, destroying the sparse structure in BA [10]. In the process of marginalization, the DSO maintains a priori information between frames (see EnergyFunctional::HM and bM) and uses this information in the solution of the BA.

2.5 Loop Closing

References

- [1] Sean L Bowman, Nikolay Atanasov, Kostas Daniilidis, and George J Pappas. Probabilistic data association for semantic slam. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 1722–1729. IEEE, 2017.
- [2] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 4, 2017.
- [3] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European Conference on Computer Vision*, pages 834–849. Springer, 2014.

- [4] Christian Forster, Matia Pizzoli, and Davide Scaramuzza. Svo: Fast semi-direct monocular visual odometry. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 15–22. IEEE, 2014.