

experimental comparison of some of these algorithms on image retrieval datasets. You can also find more details on related techniques and systems in Section 6.2.3 on visual similarity search, which discusses global descriptors that represent an image with a single vector (Arandjelovic, Gronat *et al.* 2016; Radenović, Tolias, and Chum 2019; Yang, Kien Nguyen *et al.* 2019; Cao, Araujo, and Sim 2020; Ng, Balntas *et al.* 2020; Tolias, Jenicek, and Chum 2020) as alternatives to bags of local features, Section 11.2.3 on location recognition, and Section 11.4.6 on large-scale 3D reconstruction from community (internet) photos.

7.1.5 Feature tracking

An alternative to independently finding features in all candidate images and then matching them is to find a set of likely feature locations in a first image and to then *search* for their corresponding locations in subsequent images. This kind of *detect then track* approach is more widely used for video tracking applications, where the expected amount of motion and appearance deformation between adjacent frames is expected to be small.

The process of selecting good features to track is closely related to selecting good features for more general recognition applications. In practice, regions containing high gradients in both directions, i.e., which have high eigenvalues in the auto-correlation matrix (7.8), provide stable locations at which to find correspondences (Shi and Tomasi 1994).

In subsequent frames, searching for locations where the corresponding patch has low squared difference (7.1) often works well enough. However, if the images are undergoing brightness change, explicitly compensating for such variations (9.9) or using *normalized cross-correlation* (9.11) may be preferable. If the search range is large, it is also often more efficient to use a *hierarchical* search strategy, which uses matches in lower-resolution images to provide better initial guesses and hence speed up the search (Section 9.1.1). Alternatives to this strategy involve learning what the appearance of the patch being tracked should be and then searching for it in the vicinity of its predicted position (Avidan 2001; Jurie and Dhome 2002; Williams, Blake, and Cipolla 2003). These topics are all covered in more detail in Section 9.1.3.

If features are being tracked over longer image sequences, their appearance can undergo larger changes. You then have to decide whether to continue matching against the originally detected patch (feature) or to re-sample each subsequent frame at the matching location. The former strategy is prone to failure, as the original patch can undergo appearance changes such as foreshortening. The latter runs the risk of the feature drifting from its original location to some other location in the image (Shi and Tomasi 1994). (Mathematically, small misregistration errors compound to create a *Markov random walk*, which leads to larger drift over time.)

Jenicek, and Chum 2020) as alternatives to bags of local features, Section 11.2.3 on location recognition, and Section 11.4.6 on large-scale 3D reconstruction from community (internet) photos.

7.1.5 Feature tracking

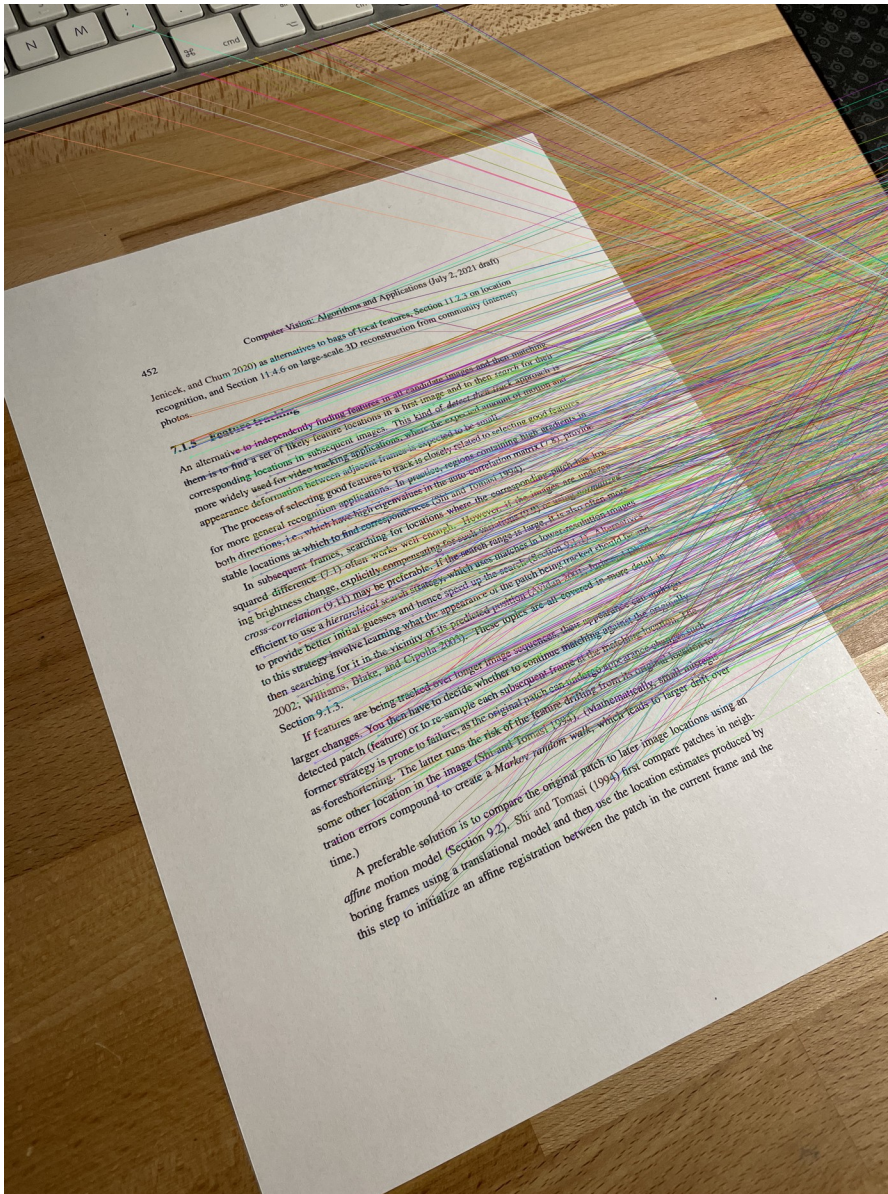
An alternative to independently finding features in all candidate images and then matching them is to find a set of likely feature locations in a first image and to then *search* for their corresponding locations in subsequent images. This kind of *detect then track* approach is more widely used for video tracking applications, where the expected amount of motion and appearance deformation between adjacent frames is expected to be small.

The process of selecting good features to track is closely related to selecting good features for more general recognition applications. In practice, regions containing high gradients in both directions, i.e., which have high eigenvalues in the auto-correlation matrix (7.8), provide stable locations at which to find correspondences (Shi and Tomasi 1994).

In subsequent frames, searching for locations where the corresponding patch has low squared difference (7.1) often works well enough. However, if the images are undergoing brightness change, explicitly compensating for such variations (9.9) or using *normalized cross-correlation* (9.11) may be preferable. If the search range is large, it is also often more efficient to use a *hierarchical* search strategy, which uses matches in lower-resolution images to provide better initial guesses and hence speed up the search (Section 9.1.1). Alternatives to this strategy involve learning what the appearance of the patch being tracked should be and then searching for it in the vicinity of its predicted position (Avidan 2001; Jurie and Dhome 2002; Williams, Blake, and Cipolla 2003). These topics are all covered in more detail in Section 9.1.3.

If features are being tracked over longer image sequences, their appearance can undergo larger changes. You then have to decide whether to continue matching against the originally detected patch (feature) or to re-sample each subsequent frame at the matching location. The former strategy is prone to failure, as the original patch can undergo appearance changes such as foreshortening. The latter runs the risk of the feature drifting from its original location to some other location in the image (Shi and Tomasi 1994). (Mathematically, small misregistration errors compound to create a *Markov random walk*, which leads to larger drift over time.)

A preferable solution is to compare the original patch to later image locations using an *affine* motion model (Section 9.2). Shi and Tomasi (1994) first compare patches in neighboring frames using a translational model and then use the location estimates produced by



experimental comparison of some of these algorithms on image retrieval datasets. You can also find more details on related techniques and systems in Section 6.2.3 on visual similarity search, which discusses global descriptors that represent an image with a single vector (Arazulov, Girona *et al.* 2016; Radhakrishna, Tofani, and Cham 2019; Yang, Kim, Nguyen *et al.* 2019; Cui, Araya, and Sim 2020; Ng, Baheti *et al.* 2020; Tofani, Jenček, and Cham 2020) as alternatives to bags of local features, Section 11.2.3 on location recognition, and Section 11.4.6 on large-scale 3D reconstruction from community (internet) photos.

7.1.5 Feature tracking

Feature tracking is independently finding features in all candidate images and then matching them to find a set of likely feature locations in a first image and to then search for their corresponding locations in subsequent images. This kind of *feature matching* approach is more widely used for video tracking applications, where the expected amount of motion and appearance deformation between adjacent frames is expected to be small.

The increased effectiveness of feature matching is closely related to selecting good features for feature matching applications. In feature matching, features containing high gradients in both directions, i.e., which have high eigenvalues in the auto-correlation matrix (Eq. 7.1), provide stable locations at which to find correspondences (Shi and Tomasi 1994).

In feature matching, searching for locations where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024). Alternatives to feature matching are feature matching, which uses matching in lower-resolution images (e.g., 64 × 64 pixels), where the number of pixels is small (e.g., 64 × 64). Alternatives to feature matching are feature matching, which uses matching in lower-resolution images (e.g., 64 × 64 pixels), where the number of pixels is small (e.g., 64 × 64).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024). Alternatives to feature matching are feature matching, which uses matching in lower-resolution images (e.g., 64 × 64 pixels), where the number of pixels is small (e.g., 64 × 64).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).

Feature matching is often used for feature matching, where the corresponding patch has low squared difference (Eq. 7.1) often comes with a cost. However, if the images are undistorted and brightness changes are small (compensating for such variations by using normalized intensity differences as in Eq. 7.1) can be made to be much less expensive. It is also often more expensive to find feature matching locations, which are often in high-resolution images (e.g., 1024 × 1024 pixels), where the number of pixels is large (e.g., 1024 × 1024).