

Sequential localisation and map-building for real-time computer vision and robotics[☆]

Andrew J. Davison^a, Nobuyuki Kita^{b,*}

^a Robotics Research Group, University of Oxford, Parks Road, Oxford OX1 3PJ, UK

^b Intelligent Systems Laboratory, National Institute of Advanced Industrial Science and Technology, AIST Tsukuba Central 2, Tsukuba 305-8568, Japan

Received 11 August 2000; received in revised form 3 April 2001

Communicated by F.C.A. Groen

Abstract

Reviewing the important problem of simultaneous localisation and map-building, we emphasise its genericity and in particular draw parallels between the often divided fields of computer vision and robot navigation. We compare sequential techniques with the batch methodologies currently prevalent in computer vision, and explain the additional challenges presented by real-time constraints which mean that there is still much work to be done in the sequential case, which when solved will lead to impressive and useful applications. In a detailed tutorial on map-building using first-order error propagation, particular attention is drawn to the roles of modelling and an active methodology. Finally, recognising the critical role of software in tackling a generic problem such as this, we announce the distribution of a proven and carefully designed open-source software framework which is intended for use in a wide range of robot and vision applications. © 2001 Elsevier Science B.V. All rights reserved.

Keywords: Navigation; Map-building; Localisation; Sequential versus batch; Active sensing

1. Introduction

Structure from motion in computer vision and simultaneous map-building and localisation for mobile robots (“SLAM”) are two views of the same problem. The situation under consideration is that of a body which moves through a static environment about which it has little or no prior knowledge, and measurements from its sensors are used to provide information about its motion and the structure of the world. This body could be a single camera, a robot with various sensors, or any other moving object able to sense its surroundings. Nevertheless, in recent years, research has taken many paths for solving this problem with a lack of acknowledgement of its general nature, with a particular divide arising between robotic and “pure vision” approaches.

Crucially, in this paper, we are interested in the *sequential* case, where map-building and localisation are able to proceed in a step-by-step fashion as movement occurs. We will contrast this with situations where the *batch*

[☆] The authors are grateful to Philip McLauchlan, Simon Julier and John Leonard for discussions.

* Corresponding author. Tel.: +81-298-61-5986; fax: +81-298-61-5971.

E-mail addresses: ajd@robots.ox.ac.uk (A.J. Davison), n.kita@aist.go.jp (N. Kita).

methods currently prevalent in computer vision (and their cousins in robot map-building) can be applied, where measurements from different time steps are used in parallel after the event. Despite renewed interest in sequential map-building from the robotics community, in computer vision, recent successful work in off-line reconstruction from image sequences has conspicuously not been accompanied by advances in real-time methods. Sequential map-building is a problem which is far from being solved, and we will look at the state of the art and its limitations.

1.1. Applications

In map-building applications where localisation or map estimates are needed quickly and successively, either to supply data to external processes in real-time or to feed back into determining future actions, *only sequential methods can be used*. For instance:

- Camera-based structure from motion methods that need to update in real-time, like in an inside-out head-tracking application, where an outward-looking camera attached to a head-mounted display user's head identifies and tracks arbitrary features in the surroundings to calculate head movement, or in live virtual studio applications, where the movement of television cameras needs to be known precisely so that live and computer-generated images can be fused to form composite output.
- Autonomous robot navigation in unknown environments, where sensor readings are used to build and update maps, and continually estimate the robot's movement through the world.

1.2. Aims of this paper

1. To review and clarify the status of the sequential map-building problem, and emphasise its genericity within robotics and computer vision.
2. In a detailed tutorial on map-building using first-order error propagation, to discuss a number of details about implementing real sequential systems and explain the approaches which our experience has led to.
3. To announce the distribution of an open-source software package for sequential localisation and map-building, designed with acknowledgement of the general nature of this class of problems and therefore readily applicable in many applications, and already proven in varied research projects [6,7,14].

2. The challenges of sequential map-building

2.1. The key point

Thinking first not of actively adding to a map, but of updating uncertain estimates of the locations of various features and that of a moving sensor platform measuring them in a sequential, real-time sense, *the amount of computation which can be carried out in each time-step is bounded by a constant*. This follows simply from thinking of implementing such a system, however fast the processor available, it can only do so much in a certain time-step.

A major implication of this is that *at a given time, we must express all our knowledge of the evolution of the system up to that time with an amount of information bounded by a constant*. The previous knowledge must be combined with any new information from the current time-step to produce updated estimates within the finite processing time available.

In the following sections, we will look at the approach that we are forced to take to fit this constraint and the difficulties of this presents, since in sequential processing compromises must be made to maintain the processing speed.

2.2. Approaching sequential map-building

There would seem to be a difference between robot map-building for localisation, where the goal is to determine the robot's motion making use of arbitrary features in the environment as landmarks, and structure from motion, where the interest is in the scene structure and not in the arbitrary path of the camera used to study it. However, it is necessary, explicitly or implicitly, to estimate both sensor motion and scene structure together if either is to be determined.

2.2.1. Batch methods

The optimal way to build maps from measurements from a moving robot or sensor is to take all the data obtained over a motion sequence and analyse it *at once*. Estimates of where the robot was at each measurement location are calculated altogether, along with the locations of the features observed, and then adjusted until they best fit all the data. This is the *batch* methodology which is used in state-of-the-art geometrical vision to recover 3D structure maps from video sequences and auto-calibrate cameras (e.g. [18,21]). In robot navigation, batch methods have a shorter history but have appeared recently under the banner *EM* [20], where maps of natural features were formed from a data set collected in an initial guided robot tour of an area; afterwards the robot could use the map during autonomous navigation. However, while batch methods can build optimal maps from previously acquired data sets, they do not offer a way to incrementally change maps in the way required by real-time applications as new data is acquired. The key to why not is that the processing effort needed to calculate an estimate for each robot or camera location on a trajectory depends on the *total number* of locations. If the robot or camera moves to a new location and we wish to combine new measurement data with an existing map, all previous estimates must be revised. This does not fit our requirement for constant processing cost for sequential applications.

A comment on uncalibrated methods, which have become closely intertwined with batch estimation in computer vision: most advanced structure from motion approaches operate by assuming that certain parameters defining the camera's operation (such as the focal length) are completely unknown, and calculations take place in versions of the world which are warped in some unknown way relative to reality via the mathematics of projective geometry. Resolution to real Euclidean estimates only happens as a final step, often after an auto-calibration procedure. It should be remembered that estimating the unknown calibration parameters of a camera in this way is somewhat of a detail when it comes to the general problem of reconstructing the world from uncertain sensor measurements. These warped worlds are of no use when we wish to incorporate other types of information; this could be data from other sensors, but most importantly we mean motion information. Motion models are inextricably tied to the real, Euclidean world — the argument may be made that things such as straight line motion are preserved under projective transformation, but physics is also about rotations and scales (such as that provided by the constant of gravitational acceleration). We argue that the best course for sequential methods is to place estimates in an Euclidean frame straight away. Of course some things may be uncertain, such as absolute scales or calibration parameters, but this uncertainty can be included explicitly in estimates and reduced as more data is acquired. We do not lose any deductive power by doing this: we only add it. As will be explained below, when the interdependence between estimates is propagated properly, we retain the ability to, for instance, estimate the ratio of the depths of two features with a high precision, even if either individually is poorly known.

2.2.2. Sequential methods

To tackle the sequential case, we need a representation of the current “state” of the system whose size does not vary over time. Then, this state can be updated in a constant time when new information arrives. Both the state and the new information will be accompanied by uncertainty, and we must take account of this when weighting the old and new data to produce updates. We are taken unavoidably into the domain of time-dependent *statistics*, whereas the optimisation approach used in batch methods permits a more lax handling of uncertainty.

Something to clarify early on is that when we talk here about a state being of constant size, we mean that for a map with a given number of features the state size does not change with *time*. The fact that the state size, and

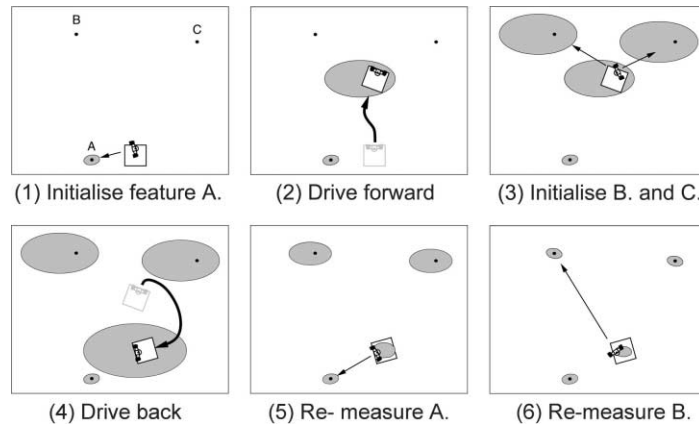


Fig. 1. Six steps in an example of sequential map-building, where a robot moving in 2D is assumed to have a fairly accurate sensor allowing it to detect the relative location of point features, and less accurate odometry for dead-reckoning motion estimation. Black points are the true locations of environmental features, and grey areas represent uncertain estimates of the feature and robot positions.

therefore processing burden, will increase as the number of features grows seems unavoidable. So to process maps in real-time, we will be limited to a finite number of features. How to deal with more features than this limit is the main challenge of sequential map-building research, and we will look at this further in Section 2.4.

Many early authors [1,8,11] took simple approaches for representing the state and its uncertainty; the locations of the moving robot in the world and features were stored and updated independently. However, if any type of long-term motion is attempted, these methods prove to be deficient: though they produce good estimates of instantaneous motion, they do not take account of the interdependence of the estimates of different quantities, and maps are seen to drift away from ground truth in a systematic way, as it can be seen in the experiments of the authors referenced above. They are not able to produce sensible estimates for long runs where previously seen features may be revisited after periods of neglect, an action that allows drifting estimates to be corrected [7].

To give a flavour of the interdependence of estimates in sequential map-building, and emphasise that it is important to estimate robot and feature positions together, steps from a simple scenario are depicted in Fig. 1. The sequence of robot behaviour here is not intended to be optimal; the point is that a map-building algorithm should be able to cope with arbitrary actions and make use of all the information it obtains.

In (1), a robot is dropped into an environment of which it has no prior knowledge. Defining a coordinate frame at this starting position, it uses a sensor to identify feature A and measure its position. The sensor is quite accurate, but there is some uncertainty in this measurement which transposes into the small grey area representing the uncertainty in the estimate of the feature's position.

The robot drives forward in (2), during this time making an estimate of its motion using dead-reckoning (for instance counting the turns of its wheels). This type of motion estimation is notoriously inaccurate and causes motion uncertainties which grow without bound over time, and this is reflected in the large uncertainty region around the robot representing its estimate of its position. In (3), the robot makes initial measurements of features B and C. Since the robot's own position estimate is uncertain at this time, its estimates of the locations of B and C have large uncertainty regions, equivalent to the robot position uncertainty plus the smaller sensor measurement uncertainty. However, although it cannot be represented in the diagram, the estimates in the locations of the robot, B and C are all *coupled* at this point. Their relative positions are quite well known; what is uncertain is the position of the group as a whole.

The robot turns and drives back to near its starting position in (4). During this motion its estimate of its own position, again updated with dead-reckoning, grows even more uncertain. In (5) though, re-measuring feature A, whose absolute location is well known, allows the robot dramatically to improve its position estimate. The important

thing to notice is that this measurement *also improves the estimate of the locations of features B and C*. Although the robot had driven farther since first measuring them, estimates of these feature positions were still partially coupled to the robot state, so improving the robot estimate also upgrades the feature estimates. The feature estimates are further improved in (6), where the robot directly re-measures feature B. This measurement, while of course improving the estimate of B, also improves C due to their interdependence (the relative locations of B and C are well known).

At this stage, all estimates are quite good and the robot has built a useful map. It is important to understand that this has happened with a small number of measurements because use has been made of the coupling between estimates.

2.3. Propagating coupled estimates

To work with coupled estimates, it is necessary to propagate not only each estimated quantity and its uncertainty, but also how this relates to the uncertainties of other estimates. Generally, a group of uncertain quantities is represented by a probability distribution in multiple dimensions, the form of which will depend on the specific agents of uncertainty in the system. Representing arbitrary probability distributions is not straightforward: one approach uses many random samples (sometimes called particles) to build up the shape, and has recently successfully been seen in vision in the form of the Condensation algorithm [12] for robust contour tracking. This approach has also been used in robot navigation for the problem of localisation using a known map [10], performing extremely well even for the difficult problem of re-localising a robot which is completely lost. However, these Monte Carlo methods are computationally expensive, and particularly are not applicable to the very high-dimensional map-building problem, since the number of particles N , and therefore the computational burden, needed to represent fairly a probability distribution in dimension d varies as

$$N \geq \frac{D_{\min}}{\alpha^d},$$

where D_{\min} and α are constants with $\alpha \ll 1$. In modern Condensation applications, the number of dimensions under consideration is limited by computing power to perhaps something less than 10 if real-time operation is desired. This is of course sufficient for estimating the location of a robot with a known map, but not when we simultaneously need to estimate map parameters.

Currently more feasible is to propagate first-order approximations to probability distributions. Each estimated parameter is accompanied by single numbers representing its variance and covariance with other parameters — a vector of parameters has a covariance matrix filled with these elements. The Kalman filter is an optimal solution to linear problems in which all noise sources are gaussian in profile, however, most map-building scenarios are not linear, so in these cases, the extended Kalman filter (EKF) provides an approximation which in general has been found to perform very well. Called “Stochastic Mapping” in its first correctly formulated application to robot map-building [19], the EKF has been implemented successfully in different scenarios by other researchers [2,4,5,7,9,15]. Its main weakness compared to the Monte Carlo methods is its inability to represent multi-modal distributions — where an estimate has two or more peak values that are most likely, with unlikely regions in between. The Monte Carlo methods gain greatly in robustness through this ability. In first-order approaches, multiple hypotheses must be considered externally and explicitly.

2.4. The problems

We consider that there are two main challenges in sequential map-building using first-order uncertainty — these are two things we have already touched on:

1. The growth of computational complexity with the number of map features.
2. Coping with mismatches (sometimes known as data association).

We will see in Section 3 that correctly propagating the interdependence of map estimates requires a covariance matrix to be updated whose number of elements grows as the square of the number of estimated parameters. Clearly, the processing requirements quickly get out of hand as the map grows.

The first approach which can be taken is to keep the number of features low using sensible map management. With today's processing power, there need be little trouble with maintaining maps with features in the 50–100 range at decent update rates. This is quite sufficient for many localisation tasks within the confined areas: the emphasis should be on using a small number of high-quality features, and on *map using* rather than wasteful map-building. For the situation in question, the question “how many features need to be measurable at any given time?” can be asked. The answer will depend on the localising power of a single feature measurement, which depends on the sensor and feature types (for instance, for a camera moving in 3D, seeing just a single point feature will not improve estimates along the several degrees of freedom which that measurement does not constrain, whereas for a robot moving in one dimension with a range sensor, measuring one feature tells it all it needs to know), as well as potential desires for redundant measurements to improve robustness (see Section 3.5). A management algorithm can then add new features to the map only in places where less than this desired number is available.

For applications where the number of features must be larger, various authors have looked at ways to relieve the complexity of large maps. One simple approach, similar to the map management above, is to delete features from the map which do not provide much information [9]: for instance, if two features lie close together, and their estimates have become closely coupled, one can be deleted without sacrificing much information. This opens up the question of which features provide the best localisation information — something also looked at in [5] with regard to active choice among candidate measurements.

Other approaches split a large map into sub-maps, within each of which fully coupled map-building goes on as normal but between which full coupling is not maintained [4,16]. While in a certain area, the robot will only observe features in the current sub-map and only update estimates for this sub-map in real-time. This is an approximation, since in truth each feature estimate is related to every other in the world map, but can be effective when the coupling between sub-maps can be represented as a single parameter: all the estimates in sub-map A have a similar dependence on all of those in sub-map B. This is likely to be the case if the robot spends long periods of time confined to each sub-map region rather than frequently moving between them. There are still many issues to settle with these approaches, such as how sub-map areas should be delineated.

McLauchlan's VSDF [17] is a powerful framework which marries sequential and batch methods, and has been used in several different vision applications. It is based on the propagation of inverse covariance matrices (called information matrices), a strategy which provides some computational advantages, and offers an efficient sequential mode, though this mode makes implicit approximations by ignoring some non-diagonal matrix elements and it is not clear how these approximations compare with other possibilities. Covariance intersection [24], a general tool for distributed estimation, has also been touted as suitable for map-building, though its generalisations may be too forgiving to produce estimates as good as other more specific methods.

A final approach [5] suggests that while maintaining a full single map, current active parts can be kept fully up-to-date while currently uninteresting regions are kept on “the back burner” — or perhaps in a hierarchy of up-to-date-ness — and proves that it is possible to bring these back up-to-date at a later stage with no loss of information. This method is certainly interesting, particularly because the hierarchy idea provides possibilities for multiple-hypothesis branching at various levels, but presents large challenges in terms of management: deciding which parts should be updated or left to simmer.

We will look in Section 3.5 that how it is possible to use redundancy and methods like RANSAC to reduce the chances of falling prey to mismatches at the measurement stage, but more generally dealing with this problem requires a multiple-hypothesis approach where estimates fork and the decision on which branch is correct is postponed until more evidence is available. There have yet to be any convincing demonstrations of how this can be incorporated into rigorous sequential map-building. This should be a major focus of future research, since there is no point in improving efficiency with methods which are still prone to instant failure at a mismatch. In our experience,

map-building of the type described in Section 3 often can survive a mismatch, though this is by luck since the method includes no model of these events.

3. Map-building with first-order uncertainty propagation: A tutorial

In the following section, we will look in detail at sequential map-building using first-order uncertainty propagation. On its own, this represents an obvious and rigorous approach to map-building, but it is also the backbone of the methods described in the previous section for improving efficiency. We will refer throughout to the moving sensor body as “the robot”, but this can apply to a single camera or other unrobot-like object. “Feature” is also a general term referring to any object the robot is capable of observing. Features can be points, lines, planes, cylinders or any other type of geometrical object.

3.1. The state vector and its covariance

Current estimates of the state of the robot and the scene features which are known about are stored in the system state vector $\hat{\mathbf{x}}$ and the uncertainty of the estimates in the covariance matrix \mathbf{P} . $\hat{\mathbf{x}}$ and \mathbf{P} will change in size dynamically as features are added to or deleted from the map. They are partitioned as follows:

$$\hat{\mathbf{x}} = \begin{pmatrix} \hat{\mathbf{x}}_v \\ \hat{\mathbf{y}}_1 \\ \hat{\mathbf{y}}_2 \\ \vdots \end{pmatrix}, \quad \mathbf{P} = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy_1} & \mathbf{P}_{xy_2} & \cdots \\ \mathbf{P}_{y_1x} & \mathbf{P}_{y_1y_1} & \mathbf{P}_{y_1y_2} & \cdots \\ \mathbf{P}_{y_2x} & \mathbf{P}_{y_2y_1} & \mathbf{P}_{y_2y_2} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

$\hat{\mathbf{x}}_v$ is the robot state estimate, and $\hat{\mathbf{y}}_i$ the estimated state of the i th feature. By the “state” of the robot and features, generally we mean a vector of all the parameters in which we are interested relating to those objects. Of course, this means their positions, which can be defined by a number of parameters depending on the geometrical type of the object and dimensionality of the map; but also, there may be other parameters which we would like the estimate, usually because they will affect future motion or measurements.

For dynamically moving objects, it is necessary to estimate higher-order motion parameters (velocity, acceleration, etc.). The number of derivatives needed depends on the expected motion (see Section 3.3). As another example, a robot may have redundant axes of movement whose status is important but which are not uniquely defined by the robot’s geometrical position. These extra parameters can also be used for *calibration constants*, which are initially only approximately known but whose accuracy will improve with the evolution of the system.

3.2. Coordinate frames and initialisation

When a robot moves in the surroundings which are initially completely unknown, the choice of a world coordinate frame is arbitrary. The only things that can be reported on are the location of the robot relative to any features detected. Indeed, one possible approach is to do away with a world coordinate frame altogether and estimate just the locations of features in a frame which is fixed to the robot; robot motions appear simply as backward feature movements. However, there is not a large computational penalty in including an explicit robot state, and more importantly in most applications of map-building, there will be some interaction with information from other sources, which could be in the form of some prior-known feature positions, or maybe simply metric way-points through which the robot is required to move — a world coordinate frame is necessary to interact with information of this kind.

If there is no prior knowledge of the environment, the coordinate frame can be defined to have its origin at the robot’s starting position, and the initial uncertainty relating to the robot’s position in \mathbf{P}_{xx} is set to zero. If there

is prior knowledge of some feature locations, this is put into the map explicitly and it is this which defines the coordinate frame. The robot's starting position relative to these features must also be input, and both robot and feature positions should be assigned as suitable initial covariance values. It is not reasonable to set both the robot and feature covariances to zero, because their relative locations can never be perfectly known; a typical initial situation would be to have several very well-known feature positions with low covariance effectively pinning down the coordinate frame, with a more uncertain robot starting location.

3.3. Motion

What happens to the estimate of the robot's position during a movement? The answer is that we should model its movement as best as possible with a motion model $f_v(\mathbf{x}_v, \mathbf{u})$, and add to the system covariance to account for our uncertainty in this motion estimate (the process noise \mathbf{Q}).

In batch structure from motion, there is typically no motion modelling. The assumption made is that at each new camera position, there is no prior location knowledge; that is to say, there is infinite uncertainty (though there may be constraints on some movement dimensions in certain scenarios). In the quasi-static case where these methods are applied to this is sensible. However, when working in the time-domain, there is always extra information to be had by modelling motion. This model may be very simple or vague, but the best thing to do is to set it up as honestly as possible and make use of it. Quoting from Torr et al. [22], who in turn cites Jaynes [13]:

Some will complain that to use Bayesian methods one must introduce arbitrary priors on the parameters. However, far from being a disadvantage, this is a tremendous advantage as it forces open acknowledgement of what assumptions were used in designing the algorithm, which all too often are hidden away beneath the veneer of equations.

There are many types of motion model depending on the level of our knowledge about the system. In the case that we have knowledge of the control parameters of a robot (such as "drive forward at 1 ms^{-1} for one second with a steering angle of 5° "), which is the case for a robot which is controlling its own navigation, we can potentially make quite an accurate motion estimate and the process noise covariance will be small. However, if we want to estimate the motion of a camera strapped to an independently manoeuvring human head, we can make much less precise assumptions: for instance, that the head will keep moving more or less at its current speed, or maybe that it is slightly more likely to slow down than speed up, given that the person is probably moving within a confined space (see the auto-regressive models of [12]).

Process noise accounts for the things that we do not attempt to model. There is of course no such thing as random noise in (classical!) physics: bodies move deterministically. It is just that we cannot know all the details of what is happening to them, though in theory we could model everything (slipping wheels; human joints, muscles, chemical reactions!). Models have to stop somewhere: the rest of what is going on we call process noise, and estimate the size of its effect.

The robot's motion is discretised into finite steps, with an incrementing label k affixed to each. There is no need for the length of these steps to be equal in time, though this will often be the case. (One point we would like to highlight is that many navigation researchers have used unnecessarily simple motion models for their mobile robots; e.g. [8], where a model for car-like motion is used which is an approximation for small Δt : in this case, it is quite straightforward to construct a motion model which does not require short time steps [5].) In a Jacobian calculation, we change the state and covariance as follows:

$$\hat{\mathbf{x}}(k+1|k) = \begin{pmatrix} f_v(\mathbf{x}_v(k|k), \mathbf{u}(k)) \\ \hat{\mathbf{y}}_1(k|k) \\ \hat{\mathbf{y}}_2(k|k) \\ \vdots \end{pmatrix}, \quad (1)$$

$$\mathbf{P}(k+1|k) = \begin{bmatrix} \frac{\partial \mathbf{f}_v}{\partial \mathbf{x}_v} \mathbf{P}_{xx}(k|k) \frac{\partial \mathbf{f}_v^T}{\partial \mathbf{x}_v} + \mathbf{Q}(k) & \frac{\partial \mathbf{f}_v}{\partial \mathbf{x}_v} \mathbf{P}_{xy_1}(k|k) & \frac{\partial \mathbf{f}_v}{\partial \mathbf{x}_v} \mathbf{P}_{xy_2}(k|k) & \cdots \\ \mathbf{P}_{y_1x}(k|k) \frac{\partial \mathbf{f}_v^T}{\partial \mathbf{x}_v} & \mathbf{P}_{y_1y_1}(k|k) & \mathbf{P}_{y_1y_2}(k|k) & \cdots \\ \mathbf{P}_{y_2x}(k|k) \frac{\partial \mathbf{f}_v^T}{\partial \mathbf{x}_v} & \mathbf{P}_{y_2y_1}(k|k) & \mathbf{P}_{y_2y_2}(k|k) & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}. \quad (2)$$

3.4. Measurements: Selection, prediction and searching

The way to measure a particular feature i is determined by its feature measurement model $\mathbf{h}_i(\mathbf{x}_v, \mathbf{y}_i)$, and the measurement noise \mathbf{R} . Analogous to process noise, measurement noise takes account of the things which we do not model in the feature measurement model. Whenever we wish to measure a particular feature, the value of the measurement can be predicted by substituting current estimates $\hat{\mathbf{x}}_v$ and $\hat{\mathbf{y}}_i$ into the expression for \mathbf{h}_i . From the predicted value of a measurement, we can calculate, based on the knowledge about the particular feature type and saved information, on what the first initialisation measurement of this feature was, whether it is worth trying to measure it. For instance, when measuring point features visually with correlation, there is little chance of a successful match if the current viewpoint is far from the original viewpoint. In this way, regions of *measurability* can be defined for each feature, and aid robustness by only allowing match attempts from positions where the chances are good.

The innovation covariance \mathbf{S}_i is how much the actual measurement \mathbf{z}_i is expected to deviate from this prediction:

$$\mathbf{S}_i = \frac{\partial \mathbf{h}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xx} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{x}_v} + \frac{\partial \mathbf{h}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xy_i} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{y}_i} + \frac{\partial \mathbf{h}_i}{\partial \mathbf{y}_i} \mathbf{P}_{y_i x} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{x}_v} + \frac{\partial \mathbf{h}_i}{\partial \mathbf{y}_i} \mathbf{P}_{y_i y_i} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{y}_i} + \mathbf{R}. \quad (3)$$

Calculating \mathbf{S}_i before making measurements allows us to form a *search region* in measurement space for each feature at a chosen number of standard deviations. This is a large advantage because it allows the adoption of an *active* approach: we need only direct searching attention to this area, maximising computational resources and minimising the chance of obtaining a mismatch.

It is also possible to make decisions on which of several potentially measurable features to observe as a priority based on \mathbf{S}_i : if the measurement cost for each candidate is similar, it is favourable to make a measurement where \mathbf{S}_i has a high value because there is most information to be gained here. There is no point in making a measurement when the result is predictable [7]. If measurements are continually chosen like this, the uncertainty in any particular part of the map can be stopped from getting out of control, a situation which would lead to large search regions and the high possibility of mismatches, and to the potential breaking of the linearisation approximations of the EKF.

3.5. Updating after a measurement

After attempting measurements, those that were successful are used to update the state estimates. How do we know which were successful? Clearly there are some cases where failures are apparent, by matching scores below a given threshold for instance. However, in other cases we would not know: something has been found within the innovation covariance-bounded search region, but is it the feature we were looking for or just something else that looks the same?

One way to pick the good measurements from the bad is make a lot of measurements at the same time, and then look for sets among them which agree with each other: these are likely to be the correct matches, since no correlation is expected amongst the failures — this algorithm, RANSAC, is commonly used to lend robustness to batch methods [23]. To use RANSAC here, we try to update with randomly selected subsets of the measurements,

and look for updated robot position estimates which agree. The bad matches can then be marked as such and the update performed with just the good ones.

To update the map based on a set of measurements \mathbf{z}_i , we perform EKF updates as below. Because the measurements are independent, these updates can be done in sequence rather than stacking all the measurements into one large vector and doing everything at once (this is computationally beneficial because smaller \mathbf{S} matrices are inverted). Note further that if a particular \mathbf{z}_i has diagonal measurement noise \mathbf{R} , we can further subdivide to the individual measurement parameters for the sequential updates.

For each independent measurement \mathbf{h}_i , the Kalman gain \mathbf{W} can be calculated and the state updated as follows:

$$\mathbf{W} = \mathbf{P} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{x}} \mathbf{S}_i^{-1} = \begin{pmatrix} \mathbf{P}_{xx} \\ \mathbf{P}_{y_1x} \\ \mathbf{P}_{y_2x} \\ \vdots \end{pmatrix} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{x}_v} \mathbf{S}_i^{-1} + \begin{pmatrix} \mathbf{P}_{xy_i} \\ \mathbf{P}_{y_1y_i} \\ \mathbf{P}_{y_2y_i} \\ \vdots \end{pmatrix} \frac{\partial \mathbf{h}_i^T}{\partial \mathbf{y}_i} \mathbf{S}_i^{-1}, \quad (4)$$

$$\hat{\mathbf{x}}_{\text{new}} = \hat{\mathbf{x}}_{\text{old}} + \mathbf{W}(\mathbf{z}_i - \mathbf{h}_i), \quad (5)$$

$$\mathbf{P}_{\text{new}} = \mathbf{P}_{\text{old}} - \mathbf{W} \mathbf{S}_i \mathbf{W}^T. \quad (6)$$

3.6. Initialising a new feature

When an unknown feature is first observed, a measurement \mathbf{z}_i is obtained from its position relative to the robot. If the measurement function $\mathbf{h}_i(\mathbf{x}_v, \mathbf{y}_i)$ is invertible to $\mathbf{y}_i(\mathbf{x}_v, \mathbf{y}_i)$, we can initialise the feature state as follows (assuming here that two features are known and the new one becomes the third):

$$\mathbf{x}_{\text{new}} = \begin{pmatrix} \mathbf{x}_v \\ \mathbf{y}_1 \\ \mathbf{y}_2 \\ \mathbf{y}_i \end{pmatrix}, \quad (7)$$

$$\mathbf{P}_{\text{new}} = \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xy_1} & \mathbf{P}_{xy_2} & \mathbf{P}_{xx} \frac{\partial \mathbf{y}_i^T}{\partial \mathbf{x}_v} \\ \mathbf{P}_{y_1x} & \mathbf{P}_{y_1y_1} & \mathbf{P}_{y_1y_2} & \mathbf{P}_{y_1x} \frac{\partial \mathbf{y}_i^T}{\partial \mathbf{x}_v} \\ \mathbf{P}_{y_2x} & \mathbf{P}_{y_2y_1} & \mathbf{P}_{y_2y_2} & \mathbf{P}_{y_2x} \frac{\partial \mathbf{y}_i^T}{\partial \mathbf{x}_v} \\ \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xx} & \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xy_1} & \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xy_2} & \frac{\partial \mathbf{y}_i}{\partial \mathbf{x}_v} \mathbf{P}_{xx} \frac{\partial \mathbf{y}_i^T}{\partial \mathbf{x}_v} + \frac{\partial \mathbf{y}_i}{\partial \mathbf{h}_i} \mathbf{R} \frac{\partial \mathbf{y}_i^T}{\partial \mathbf{h}_i} \end{bmatrix}. \quad (8)$$

It should be noted that some bias is introduced into the map in initialising features in this way if (as is usual) the measurement process is non-linear.

If \mathbf{h}_i is not invertible, it means that a single measurement does not give enough information to pinpoint the feature location (for instance, a single view of a point feature from a single camera only defines a ray on which it lies). The approach that must be followed here is to initialise it into the map as a *partially initialised feature* \mathbf{y}_{pi} , with a different geometrical type (e.g. a line feature to represent the ray we know a point must lie on), and wait until another measurement from a different viewpoint allows resolution. At this stage, a special second initialisation function $\mathbf{y}_i(\mathbf{x}_v, \mathbf{y}_{\text{pi}}, \mathbf{z}_i)$ allows the actual state \mathbf{y}_i to be determined from the partially initialised state and new measurement (feature types which require more than two steps for initialisation are also possible).

Once initialised, a feature has exactly the same status in the map as those whose positions may have been given as prior knowledge.

3.7. Deleting a feature

Deleting a feature from the state vector and covariance matrix is a simple case of removing the rows and columns which contain it. An example in a system where the second of three known features is deleted would be

$$\begin{pmatrix} x_v \\ y_1 \\ y_2 \\ y_3 \end{pmatrix} \rightarrow \begin{pmatrix} x_v \\ y_1 \\ y_3 \end{pmatrix}, \quad \begin{bmatrix} P_{xx} & P_{xy_1} & P_{xy_2} & P_{xy_3} \\ P_{y_1x} & P_{y_1y_1} & P_{y_1y_2} & P_{y_1y_3} \\ P_{y_2x} & P_{y_2y_1} & P_{y_2y_2} & P_{y_2y_3} \\ P_{y_3x} & P_{y_3y_1} & P_{y_3y_2} & P_{y_3y_3} \end{bmatrix} \rightarrow \begin{bmatrix} P_{xx} & P_{xy_1} & P_{xy_3} \\ P_{y_1x} & P_{y_1y_1} & P_{y_1y_3} \\ P_{y_3x} & P_{y_3y_1} & P_{y_3y_3} \end{bmatrix}. \quad (9)$$

In automated map-maintenance, features can be deleted if a large proportion of attempted measurements fail when the feature is expected to be measurable. This could be due to features not fitting the assumptions of their model (an assumed point feature which in fact contains regions of different depths and therefore appears very different from a new viewpoint for instance), or possibly occlusion — leading to the survival of features which do not suffer these fates.

4. Software and implementations

Acknowledgement of the generic properties of the sequential map-building problem and experience with different robot systems has led to the evolution of our original application-specific software into a general framework called *Scene*, efficiently implemented in C++ and designed with orthogonal axes of flexibility in mind:

1. Use in many different application domains; from multiple robots navigating in 1D, 2D or 3D with arbitrary sensing capabilities, to single camera.
2. Implementation of different mapping algorithms and approaches to dealing with the complexity of sequential map-building.

Scene is now available with full source code (under the GNU Lesser General Public License) at <http://www.robots.ox.ac.uk/~ajd/Scene/>. The distribution package includes interactive simulations precompiled for Linux which allow immediate hands-on experience of sequential map-building in several real and simplified problem domains, the additional tools which turn these into systems operating with real hardware, and substantial documentation.

To give an impression of how the general framework can be applied to various systems, details of some current and planned implementations, differentiated by motion and feature measurement models plugged in as modules, are given in Table 1. The simplest is a 1D test-bed, which is very useful for looking at what happens to robot and feature covariances in various situations and under different algorithms. Our main work to date has been on robot navigation using active vision [5–7], using mobile platforms which move in 2D and are equipped with steerable camera platforms which make measurements of point features in 3D with stereo vision. As detailed in [6], the software's motion model formulation is flexible enough to permit cooperative position estimation by cooperating robots, where one has stereo active vision and the other is blind, navigating primarily by odometry.

Current PC computers are powerful enough to perform correlation searches for many features at video frame rate. Our current goal is to apply the *Scene* library to real-time camera position tracking using just inside-out image measurements, potentially the “killer app” of sequential localisation and map-building, which would be useful in applications such as inside-out head-tracking or the real-time virtual studio. The first demonstrations of this type have just started to appear [3]. To be successful, it will be necessary to make use of many of the details explained in this paper. For instance, RANSAC or similar must be used to detect failed matches fast because there will not be

Table 1
Specifications for various implementations

	1D test-bed	Nomad robot with stereo active vision	Camera position tracking
World dimension	1	3	3
Position dimension	1	2	3
Motion model	Velocity only	Steer and drive	General motion in 3D
State size	1	5	12
Control size	1	3	0
Feature measurement model	Point range measurement	3D stereo point measurement	Single image point measurement
Feature dimension	1	3	3
Measurement size	1	4	2

enough processing time available to propagate multiple hypotheses. A full 3D motion model estimating velocities as well as positions must be used, and finally it will be necessary to use partially initialised feature representations to bootstrap features in 3D from multiple views.

5. Conclusions

This paper has reviewed the current state of research into sequential localisation and map-building, in particular drawing parallels between approaches in the fields of mobile robotics and visual reconstruction, and has discussed the future direction of research in these domains as they are inevitably drawn closer together.

A detailed tutorial on techniques based on first-order uncertainty propagation has been presented, drawing attention to often neglected but important aspects such as modelling and active measurements.

The discussion in this paper is backed up with a proven, freely available open-source software library for sequential localisation and map-building which is expected to be of benefit to researchers in many domains.

References

- [1] P.A. Beardsley, I.D. Reid, A. Zisserman, D.W. Murray, Active visual navigation using non-metric structure, in: *Proceedings of the Fifth International Conference on Computer Vision*, Boston, MA, IEEE Computer Society Press, Silver Spring, MD, 1995, pp. 58–65.
- [2] J.A. Castellanos, Mobile robot localization and map-building: a multisensor fusion approach, Ph.D. Thesis, Universidad de Zaragoza, Spain, 1998.
- [3] A. Chiuso, P. Favaro, H. Jin, S. Soatto, “MFM”: 3-D motion from 2-D motion causally integrated over time, in: *Proceedings of the Sixth European Conference on Computer Vision*, Dublin, 2000.
- [4] K.S. Chong, L. Kleeman, Feature-based mapping in real, large scale environments using an ultrasonic array, *International Journal of Robotics Research* 18 (2) (1999) 3–19.
- [5] A.J. Davison, Mobile robot navigation using active vision, Ph.D. Thesis, University of Oxford, 1998. <http://www.robots.ox.ac.uk/~ajd/>.
- [6] A.J. Davison, N. Kita, Active visual localisation for cooperating inspection robots, in: *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems*, Takamatsu, Japan, 2000.
- [7] A.J. Davison, D.W. Murray, Mobile robot localisation using active vision, in: *Proceedings of the Fifth European Conference on Computer Vision*, Freiburg, Germany, 1998, pp. 809–825.
- [8] H.F. Durrant-Whyte, Where am I? A tutorial on mobile vehicle localization, *Industrial Robot* 21 (2) (1994) 11–16.
- [9] H.F. Durrant-Whyte, M.W.M.G. Dissanayake, P.W. Gibbens, Toward deployments of large scale simultaneous localisation and map building (slam) systems, in: *Proceedings of the Ninth International Symposium of Robotics Research*, Snowbird, UT, 1999, pp. 121–127.
- [10] D. Fox, W. Burgard, H. Kruppa, S. Thrun, Efficient multi-robot localization based on Monte Carlo approximation, in: *Proceedings of the Ninth International Symposium of Robotics Research*, Snowbird, UT, 1999, pp. 113–120.
- [11] C.G. Harris, Geometry from visual motion, in: A. Blake, A. Yuille (Eds.), *Active Vision*, MIT Press, Cambridge, MA, 1992.
- [12] M. Isard, A. Blake, Contour tracking by stochastic propagation of conditional density, in: *Proceedings of the Fourth European Conference on Computer Vision*, Cambridge, 1996, pp. 343–356.
- [13] E.T. Jaynes, Probability theory as extended logic, Technical Report, Washington University, St. Louis, MO, 1999. <http://bayes.wustl.edu/>.

- [14] J.G.H. Knight, A.J. Davison, I.D. Reid, Constant time slam using postponement, in: *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems*, 2001, to appear.
- [15] Y.D. Kwon, J.S. Lee, A stochastic map-building method for mobile robot using 2-D laser range finder, *Autonomous Robots* 7 (1999) 187–200.
- [16] J.J. Leonard, H.J.S. Feder, A computationally efficient method for large-scale concurrent mapping and localization, in: *Proceedings of the Ninth International Symposium on Robotics Research*, Snowbird, UT, 1999.
- [17] P.F. McLauchlan, D.W. Murray, A unifying framework for structure and motion recovery from image sequences, in: *Proceedings of the Fifth International Conference on Computer Vision*, Boston, MA, IEEE Computer Society Press, Silver Spring, MD, 1995.
- [18] M. Pollefeys, R. Koch, L. Van Gool, Self-calibration and metric reconstruction in spite of varying and unknown internal camera parameters, in: *Proceedings of the Sixth International Conference on Computer Vision*, Bombay, 1998, pp. 90–96.
- [19] R. Smith, M. Self, P. Cheeseman, A stochastic map for uncertain spatial relationships, in: *Proceedings of the Fourth International Symposium on Robotics Research*, 1987.
- [20] S. Thrun, D. Fox, W. Burgard, A probabilistic approach to concurrent mapping and localization for mobile robots, *Machine Learning* 31 (1998).
- [21] P.H.S. Torr, A.W. Fitzgibbon, A. Zisserman, Maintaining multiple motion model hypotheses over many views to recover matching and structure, in: *Proceedings of the Sixth International Conference on Computer Vision*, Bombay, 1998, pp. 485–491.
- [22] P.H.S. Torr, R. Szeliski, P. Anandan, An integrated bayesian approach to layer extraction from image sequences, in: *Proceedings of the Seventh International Conference on Computer Vision*, Kerkyra, 1999, pp. 983–990.
- [23] P.H.S. Torr, A. Zisserman, Robust computation and parameterization of multiple view relations, in: *Proceedings of the Sixth International Conference on Computer Vision*, Bombay, 1998, pp. 727–732.
- [24] J.K. Uhlmann, S.J. Julier, M. Csorba, Nondivergent simultaneous map-building and localisation using covariance intersection, in: *Proceedings of AeroSense: The 11th International Symposium on Aerospace/Defense Sensing, Simulation and Controls*, SPIE, Orlando, FL, 1997, *Navigation and Control Technologies for Unmanned Systems II*.



Andrew J. Davison reads Physics at the University of Oxford, gained his B.A. (1st class honours) in 1994. His doctoral research was in the area of robot navigation using active vision. On completing his D.Phil. in 1998, Dr. Davison was awarded a European Union Science and Technology Fellowship and worked for 2 years at the Japanese Government's Electrotechnical Laboratory in Tsukuba, Japan, continuing research into robot navigation. He returned to the UK in 2000, where he is currently once again working at the University of Oxford, now on visual tracking.



Nobuyuki Kita was born in 1957. He received his B.Sc. degree in Control Engineering at Osaka University in 1979, and his M.Sc. degree from Osaka University in 1981 in computer vision research. In 1981, he joined Electrotechnical Laboratory (ETL) and studied 2D/3D image understanding. After moving to the Robotics Group of ETL in 1993, he started the active vision research and was awarded the paper prize from the Robotics Society of Japan in 1995 for a paper about binocular gaze control method. His current main interest is attention control. He is now a Senior Researcher of National Institute of Advanced Industrial Science and Technology (AIST), which is the new institute following ETL and leading a project about robotic inspection of nuclear power plant.