

MonoSLAM: Real-Time Single Camera SLAM

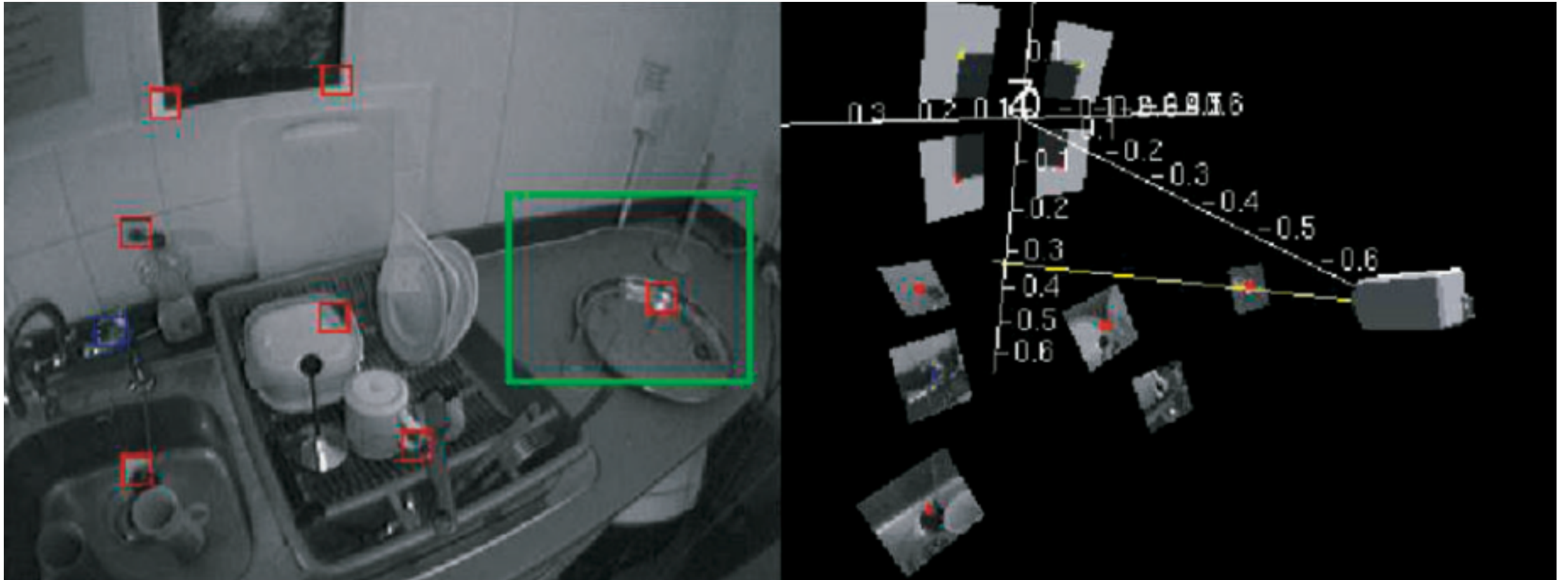
Andrew J. Davison, Ian D. Reid, *Member, IEEE*, Nicholas D. Molton, and
Olivier Stasse, *Member, IEEE*

Abstract—We present a real-time algorithm which can recover the 3D trajectory of a monocular camera, moving rapidly through a previously unknown scene. Our system, which we dub *MonoSLAM*, is the first successful application of the SLAM methodology from mobile robotics to the “pure vision” domain of a single uncontrolled camera, achieving real time but drift-free performance inaccessible to Structure from Motion approaches. The core of the approach is the online creation of a sparse but persistent map of natural landmarks within a probabilistic framework. Our key novel contributions include an *active* approach to mapping and measurement, the use of a general motion model for smooth camera movement, and solutions for monocular feature initialization and feature orientation estimation. Together, these add up to an extremely efficient and robust algorithm which runs at 30 Hz with standard PC and camera hardware. This work extends the range of robotic systems in which SLAM can be usefully applied, but also opens up new areas. We present applications of *MonoSLAM* to real-time 3D localization and mapping for a high-performance full-size humanoid robot and live augmented reality with a hand-held camera.

Index Terms—Autonomous vehicles, 3D/stereo scene analysis, tracking.



Simultaneous Localization and Mapping



Given a **single camera** feed,
estimate the 3D **position of the camera** and
the 3D **positions of all landmark** points in the world

Real-Time Camera Tracking in Unknown Scenes

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

What is the state representation?

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

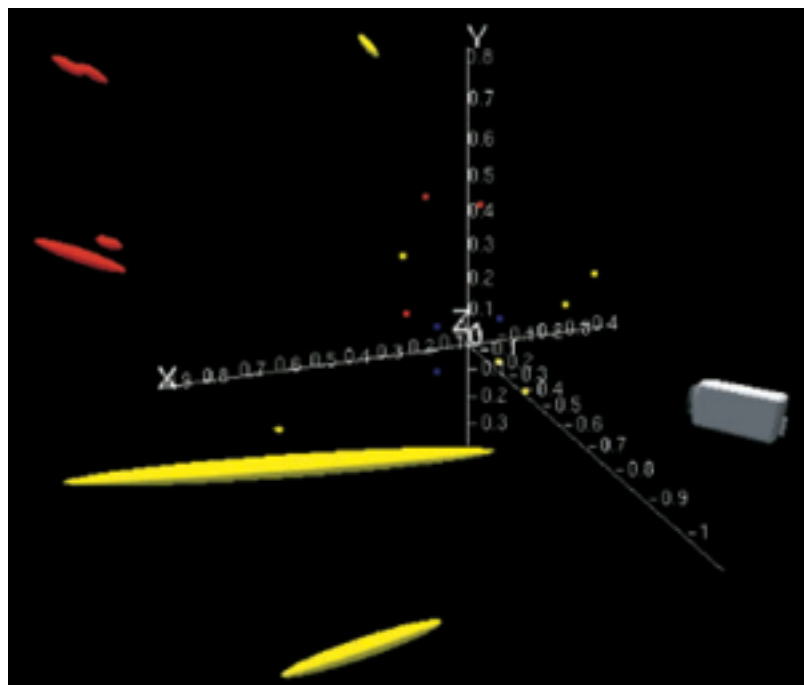
Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

What is the camera (robot) state?

What are the dimensions?

$$\mathbf{x}_c = \begin{bmatrix} \mathbf{r} \\ \mathbf{q} \\ \mathbf{v} \\ \boldsymbol{\omega} \end{bmatrix} \begin{array}{l} \text{position} \\ \text{rotation (quaternion)} \\ \text{velocity} \\ \text{angular velocity} \end{array}$$



13 total

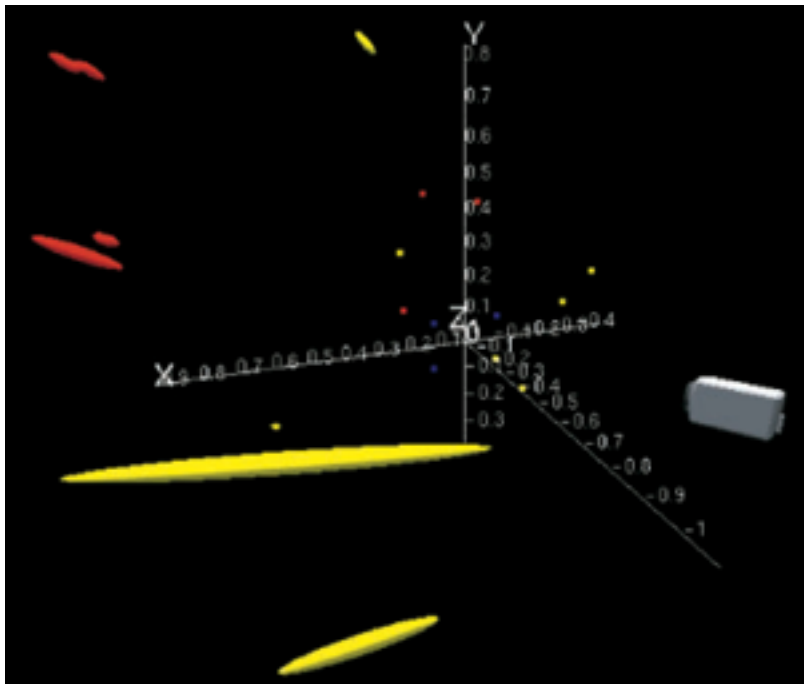
What is the camera (robot) state?

What are the dimensions?

$$\mathbf{x}_c = \begin{bmatrix} \mathbf{r} \\ \mathbf{q} \\ \mathbf{v} \\ \boldsymbol{\omega} \end{bmatrix}$$

position	3
rotation (quaternion)	4
velocity	3
angular velocity	3

13 total



What is the world (robot+environment) state?

What are the dimensions?

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_c \\ \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}$$

state of the camera

location of feature 1

location of feature 2

location of feature N

13+3N total

What is the world (robot+environment) state?

What are the dimensions?

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_c \\ \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}$$

state of the camera	13
location of feature 1	3
location of feature 2	3
location of feature N	3

13+3N total

What is the covariance (uncertainty) of the world state?

$$\Sigma = \begin{bmatrix} \Sigma_{\mathbf{x}_c \mathbf{x}_c} & \Sigma_{\mathbf{x}_c \mathbf{y}_1} & \cdots & \Sigma_{\mathbf{x}_c \mathbf{y}_N} \\ \Sigma_{\mathbf{y}_1 \mathbf{x}_c} & \Sigma_{\mathbf{y}_1 \mathbf{y}_1} & \cdots & \Sigma_{\mathbf{y}_1 \mathbf{y}_N} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{\mathbf{y}_N \mathbf{x}_c} & \Sigma_{\mathbf{y}_N \mathbf{y}_1} & \cdots & \Sigma_{\mathbf{y}_N \mathbf{y}_N} \end{bmatrix}$$

What are the dimensions?

(13+3N) x (13+3N)

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

What are the observations?

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

What are the observations?

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

Observations are...



detected visual features of landmark points.

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

What does the prediction step look like?

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

What is the motion model? $P(\mathbf{x}_t | \mathbf{x}_{t-1})$

What is the form of the belief? $P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$

What is the motion model? $P(\mathbf{x}_t | \mathbf{x}_{t-1})$

Landmarks:
constant position
(identity matrix)

Camera:
constant velocity
(not identity matrix!)

What is the form of the belief? $P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$

What is the motion model? $P(\mathbf{x}_t | \mathbf{x}_{t-1})$

Landmarks:

constant position
(identity matrix)

Camera:

constant velocity
(not identity matrix!)

What is the form of the belief? $P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$

Gaussian!

(everything is parametrized by a mean and Gaussian)

Constant Velocity Motion Model

$$\mathbf{r}_t = \mathbf{r}_{t-1} + \mathbf{v}_{t-1} \Delta t$$
 position

$$\mathbf{q}_t = \mathbf{q}_{t-1} \times [\mathbf{q}(\omega) \Delta t]$$
 rotation (quaternion)

$$\mathbf{v}_t = \mathbf{v}_{t-1}$$
 velocity

$$\omega_t = \omega_{t-1}$$
 angular velocity

Gaussian noise uncertainty (only on velocity)

$$\mathbf{v}_t = \mathbf{v}_{t-1} + \mathbf{V}$$

$$\omega_t = \omega_{t-1} + \boldsymbol{\Omega}$$

$$\mathbf{V} \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} \sigma_v & 0 & 0 \\ 0 & \sigma_v & 0 \\ 0 & 0 & \sigma_v \end{bmatrix})$$

$$\boldsymbol{\Omega} \sim \mathcal{N}(\mathbf{0}, \begin{bmatrix} \sigma_w & 0 & 0 \\ 0 & \sigma_w & 0 \\ 0 & 0 & \sigma_w \end{bmatrix})$$

Prediction (**mean** of camera state):

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

$$\mathbf{f}_t = \begin{bmatrix} \mathbf{r}_t \\ \mathbf{q}_t \\ \mathbf{v}_t \\ \omega_t \end{bmatrix} = \begin{bmatrix} \mathbf{r}_{t-1} + \mathbf{v}_{t-1} \Delta t \\ \mathbf{q}_{t-1} + \mathbf{q}(\omega)_{t-1} \Delta t \\ \mathbf{v}_{t-1} \\ \omega_{t-1} \end{bmatrix}$$

Prediction (**covariance** of camera state):

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

$$\bar{\Sigma}_{\mathbf{x}\mathbf{x}} = \boxed{\frac{\partial \mathbf{f}_t}{\partial \mathbf{x}}} \Sigma_{\mathbf{x}\mathbf{x}} \frac{\partial \mathbf{f}_t}{\partial \mathbf{x}}^\top + \mathbf{Q}_t$$

new
covariance

change
around
new state

old
covariance

change
around
new state

system noise
(process noise)

*Where does this motion model
approximation come from?*

$$\underbrace{\frac{\partial \mathbf{f}_t}{\partial \mathbf{x}_{t-1}}}_{\text{change in camera state}} = \begin{bmatrix} \underbrace{\frac{\partial \mathbf{r}_t}{\partial \mathbf{r}_{t-1}}}_{\text{change in position}} & \frac{\partial \mathbf{q}_t}{\partial \mathbf{r}_{t-1}} & \frac{\partial \mathbf{v}_t}{\partial \mathbf{r}_{t-1}} & \frac{\partial \omega_t}{\partial \mathbf{r}_{t-1}} \\ \frac{\partial \mathbf{r}_t}{\partial \mathbf{q}_{t-1}} & \frac{\partial \mathbf{q}_t}{\partial \mathbf{q}_{t-1}} & \frac{\partial \mathbf{v}_t}{\partial \mathbf{q}_{t-1}} & \frac{\partial \omega_t}{\partial \mathbf{q}_{t-1}} \\ \frac{\partial \mathbf{r}_t}{\partial \mathbf{v}_{t-1}} & \frac{\partial \mathbf{q}_t}{\partial \mathbf{v}_{t-1}} & \frac{\partial \mathbf{v}_t}{\partial \mathbf{v}_{t-1}} & \frac{\partial \omega_t}{\partial \mathbf{v}_{t-1}} \\ \frac{\partial \mathbf{r}_t}{\partial \omega_{t-1}} & \frac{\partial \mathbf{q}_t}{\partial \omega_{t-1}} & \frac{\partial \mathbf{v}_t}{\partial \omega_{t-1}} & \frac{\partial \omega_t}{\partial \omega_{t-1}} \end{bmatrix}$$

What are the dimensions?

Skipping over many details...

$$\frac{\partial \mathbf{f}_t}{\partial \mathbf{x}_{t-1}} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{I}\Delta t & \mathbf{0} \\ \mathbf{0} & \frac{\partial \mathbf{q}_t}{\partial \mathbf{q}_{t-1}} & \mathbf{0} & \frac{\partial \omega_t}{\partial \mathbf{q}_{t-1}} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} \end{bmatrix}$$

Prediction (**covariance** of camera state):

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

$$\bar{\Sigma}_{\mathbf{x}\mathbf{x}} = \boxed{\frac{\partial \mathbf{f}_t}{\partial \mathbf{x}}} \Sigma_{\mathbf{x}\mathbf{x}} \frac{\partial \mathbf{f}_t}{\partial \mathbf{x}}^\top + \mathbf{Q}_t$$

new
covariance

change
around
new state

old
covariance

change
around
new state

system noise
(process noise)

Bit of a pain to compute this term...

We just covered the **prediction** step for the camera state

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

$$\mathbf{f}_t = \begin{bmatrix} \mathbf{r}_t \\ \mathbf{q}_t \\ \mathbf{v}_t \\ \omega_t \end{bmatrix} = \begin{bmatrix} \mathbf{r}_{t-1} + \mathbf{v}_{t-1} \\ \mathbf{q}_{t-1} + \mathbf{q}(\omega)_{t-1} \\ \mathbf{v}_{t-1} \\ \omega_{t-1} \end{bmatrix}$$

$$\bar{\Sigma}_{\mathbf{x}\mathbf{x}} = \frac{\partial \mathbf{f}_t}{\partial \mathbf{x}} \Sigma_{\mathbf{x}\mathbf{x}} \frac{\partial \mathbf{f}_t}{\partial \mathbf{x}}^\top + \mathbf{Q}_t$$

Now we need to do the **update** step!

General Filtering Equations

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) \propto P(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Prediction:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t-1}) = \int_{\mathbf{x}_{t-1}} P(\mathbf{x}_t | \mathbf{x}_{t-1}) P(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}) d\mathbf{x}_{t-1}$$

Update:

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

Belief state State observation Predicted State

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$



What are the observations?

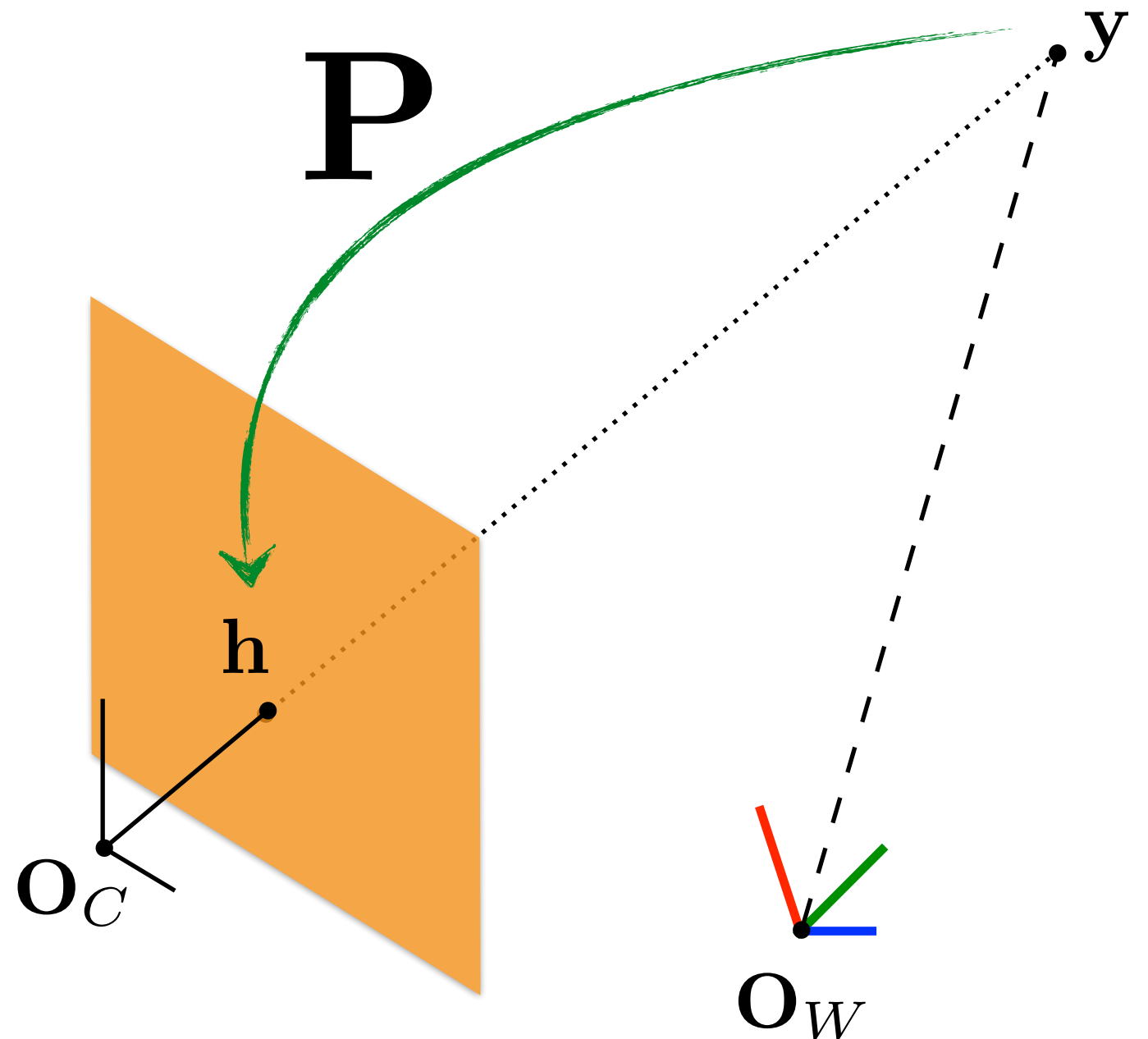


2D projections of 3D landmarks

Recall, the state includes the 3D location of landmarks

What is the projection from 3D point to 2D image point?

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_c \\ \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}$$




Observation Model

$$P(z_t | x_t)$$

If you know the 3D location of a landmark, what is the 2D projection?

**Non-linear
observation model**


$$\mathbf{h} \sim \mathbf{P} \mathbf{y}$$

2D Image Point Camera matrix 3D World Point

$$\mathbf{P} = \mathbf{K} [\mathbf{R} | \mathbf{T}]$$

*What do we know
about **P**?*

How do we make the observation model linear?

$$\begin{array}{c}
 \mathbf{H} \\
 (2n \times 13)
 \end{array}
 =
 \frac{\partial \mathbf{h}}{\partial \mathbf{x}}$$

n: number of visible points

I will spare you the pain of deriving the partial derivative...

$$P(\mathbf{x}_t | \mathbf{z}_{1:t}) = P(\mathbf{z}_t | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{z}_{1:t-1})$$

Update step (mean):

$$\underset{\text{Updated state}}{\mathbf{x}_t} = \underset{\text{Predicted state}}{\mathbf{x}_t} + \overset{\text{Kalman gain}}{\mathbf{K}_t} \left(\underset{\text{Matched 2D features}}{\mathbf{z}_t} - \underset{\text{2D projection of 3D point}}{\mathbf{h}(\mathbf{y}; \mathbf{x}_t)} \right)$$

Update step (covariance):

$$\underset{\text{Covariance (updated)}}{\Sigma_t} = \left(\overset{\text{Identity}}{\mathbf{I}} - \overset{\text{Kalman gain}}{\mathbf{K}_t} \underset{\text{Jacobian}}{\mathbf{H}_t} \right) \underset{\text{Covariance (predicted)}}{\Sigma_t}$$

Kintinuous: Spatially Extended Kinect Fusion

Thomas Whelan, John McDonald

National University of Ireland Maynooth, Ireland

Michael Kaess, Maurice Fallon, Hordur Johannsson,
John J. Leonard

Computer Science and Artificial Intelligence
Laboratory, MIT, USA

