# Assimilation of sea-surface temperature into a hydrodynamic model of Port Phillip Bay, Australia

Matthew R.J. Turner<sup>1</sup>, Jeffrey P. Walker<sup>1</sup>, Peter R. Oke<sup>2</sup> and Rodger B. Grayson<sup>1</sup>

<sup>1</sup>Department of Civil and Environmental Engineering, The University of Melbourne <sup>2</sup> CSIRO Marine Research, Hobart, Tasmania

# Abstract

Hydrodynamic models can predict states of interest to the coastal engineer, however, due to uncertainties in the model physics, model parameters, initial conditions, and model forcing data, large errors in prediction often result. To counter this, an ensemble sequential data assimilation scheme has been applied to the Model for Estuaries and Coastal Oceans (MECO), to constrain model predicted water temperature with remotely sensed sea-surface temperature observations. This paper describes a series of synthetic twin experiments that contrast two ensemble sequential data assimilation schemes, the Ensemble Kalman Filter (EnKF) and the Ensemble Square Root Filter (EnSRF) in both one and three dimensional forms. The experiments show that the assimilation greatly improves the model prediction. The three dimensional form outperforms the one dimensional form, and that the EnSRF outperformed the EnKF significantly in the one dimensional form but only marginally in the three dimensional form.

# 1 Introduction

Data assimilation is a statistical technique, which combines a model forecast and an observation to estimate the true state of the phenomenon being predicted. This paper presents an introduction to the concepts of sequential data assimilation and compares two well-known techniques from the literature in a coastal application. Synthetic surface temperature fields (SST) are assimilated into a coastal hydrodynamic model and show significant improvement in the prediction capability of the model.

The benefits of using data assimilation are its ability to improve model prediction. This is of importance in short range forecasting where a prediction of a future state is desired. In a coastal setting this could be where will be the location of an algal blooms most likely be, or what concentration of suspended sediment should we expect at a given location two days hence, while in a ports setting improved prediction of water levels or wave fields have obvious benefits for the safety and reliability of shipping. Other benefits of data assimilation are that an investigation of the data assimilation analysis can point to potential model deficiencies. For instance, if the assimilation always corrects the model in a certain direction this is suggestive of a poor model parameterisation. Thus data assimilation provides feedback that aids model improvement.

Accurate prediction of water temperature is of particular importance in ecological modelling, where temperature influences growth parameters. prediction Unfortunately, accurate of water temperature is not always possible, especially in highly enclosed water bodies where atmospheric exchange is the dominant driver of water temperature. Errors in water temperature prediction are due to our poor understanding and conceptualisation of thermodynamics and hydrodynamics, as well as uncertainties in initial conditions, model parameters, and model forcing data.

Model prediction errors can be reduced by using data assimilation. Data assimilation combines the model predicted states with observations based on their relative uncertainties. The result of data assimilation is a new set of state estimates that are closer to the truth and have a lower level of uncertainty than either data set (model or observations) individually. Data assimilation techniques have been widely applied in meteorology and oceanography (Gill and Malanotte-Rizzoli, 1991), but assimilation of sea-surface temperature (SST) into coastal ocean models has received far less attention.

Applying data assimilation to coastal models has become increasingly accessible due to recent advances in computing power and the launch of new satellite observing systems. However, while many data assimilation approaches exist in the literature, it is not clear which of these are best suited to bay and estuary modelling. Therefore, as well as introducing data assimilation generally, this paper explores the characteristics of two ensemble sequential data assimilation techniques – the Ensemble Kalman Filter (EnKF) and the Ensemble Square Root Filter (EnSRF). The key difference between these two data



Figure 1: Location diagram for Port Phillip Bay, south eastern Australia.

assimilation techniques is their treatment of observations. These two techniques are contrasted in a series of synthetic twin experiments in both one dimensional (where only the vertical correlation of the water column is considered for individual model cells and the assimilation is performed cell wise through the model domain) and three dimensional (where the spatial correlation of model states is considered and the entire model domain is modified by the assimilation in a single step) form. The experiments are undertaken for Port Phillip Bay located in south eastern Australia (Figure 1).

### 2 Theoretical Background

It is well understood that both model predictions and observations of a physical state are often prone to error. For models the mathematical equations representing the phenomenon are a simplification of actual processes; computing power limits the spatial and temporal resolution and uncertainties associated with boundary and initial conditions all combine to produce uncertainties in the model prediction.

In the coastal marine setting two data types are most commonly available; point source and remotely sensed. While point measurements are usually relatively accurate at the local scale, extending these data introduces uncertainties of scale. In contrast, remotely sensed data provide good spatial coverage, but only for a shallow surface layer at an instant in time, and are sensitive to atmospheric effects and errors in the algorithm used to relate the measurements to the physical state being observed.

While an estimate of the spatial and temporal variation in water temperature can be made based on either model predictions or observations alone, both are affected by different types of uncertainty. Observations, although subject to errors may be accurate: models give temporal state estimates for the entire domain which observations can not. Sequential data assimilation combines the model predictions and observations to achieve an improved estimate of the physical state.

The algorithm for sequential data assimilation is as follows. Starting from a 'best estimate' of the physical state, a model run predicts the physical state in the future. When an observation becomes available the model is stopped. The model state at this point becomes the 'forecast', or background field. Based on the relative uncertainty in and the difference between the observation and the forecast; and the covariances between the forecast and observation errors, a correction is calculated for the model state. This correction is added to the forecast to give the 'analysis'. The model is then reinitialised using the analysis, and is run forward until another observation becomes available and the process repeated. As the name suggests, the observations are sequentially assimilated into the model.

The most well known sequential data assimilation technique is the standard Kalman Filter (Evensen 2003) given by

$$\mathbf{x}^{a} = \mathbf{x}^{f} + \mathbf{K}(\mathbf{d} - \mathbf{H}\mathbf{x}^{f}), \qquad (1)$$

where  $\mathbf{x}$  is a vector of the model predictions, with superscripts *a* and *f* denoting analysis and forecast respectively, **d** is a vector of observations and **H** is a matrix that maps the model state  $\mathbf{x}$  to the observations **d**; **K** is a weighting matrix known as the Kalman gain given by

$$\mathbf{K} = \mathbf{P}\mathbf{H}^{\mathrm{T}}(\mathbf{H}\mathbf{P}\mathbf{H}^{\mathrm{T}} + \mathbf{R})^{-1}, \qquad (2)$$

where **P** is the forecast error covariance matrix that quantifies the covariances of the uncertainties of the model state, and **R** is the observation error covariance matrix that quantifies the covariances of the uncertainties associated with the observations. The influence of the Kalman gain on the analysis can be seen by considering an example where **P** and **R** are scalars. If the uncertainty associated with the model is less than the uncertainty associated with the model is placed on the model forecast. Conversely, if the observations are more certain than the model forecast, **R** << **P**, **K** approaches unity and more reliance is placed on the observations.

An advantage of the Kalman Filter over other sequential data assimilation techniques, such as direct insertion, is that the statistical relationship between state elements enables the filter to update not just those state elements that are observed, but also other unobserved state elements that may be different variables and at different locations. For instance, satellites typically measure SST, however using the Kalman Filter SST observations can be used to modify all other state elements, including temperature at depth, salinity, sea-level and currents.

Because the Kalman Filter is linear it does not deal satisfactorily with highly nonlinear models. In an attempt to overcome this limitation, Evensen (2003) introduced the Ensemble Kalman Filter (EnKF), whereby covariance error statistics were obtained through the use of an ensemble of model forecasts. Thus rather than one model run being propagated through time an ensemble of model runs are made. Each run starts from a slightly different position, reflecting the uncertainty associated with the model initial conditions and each model is forced by slightly different forcing data reflecting the uncertainty associated with the forcing data. Model error is incorporated too. The result is that when an observation becomes available for analysis the ensemble of forecasts will have spread representing the uncertainty associated with the forecast. This is illustrated in Figure 2 where each solid dot represents



Figure 2 Graphic representation of sequential assimilation. Dots represent ensemble members, both initial/analysed values (solid) and forecast values (open). The relative spread of an ensemble indicates its uncertainty.

one ensemble member and in combination represent the uncertainty associated with a particular model state. Each ensemble member is propagated through time by the model to give a forecast (open dot) when an observation becomes available. Due to the uncertainties in the model and forcing conditions the uncertainty of the forecast has spread. The assimilation reduces the spread of the ensemble members (solid dots) indicating a reduction in uncertainty associated with the analysed state. The analysed values are used to initiate the next forecast and the process repeats itself.

In ensemble form equation (1) becomes

$$\mathbf{X}^{a} = \mathbf{X}^{f} + \mathbf{P}_{e}\mathbf{H}^{T}(\mathbf{H}\mathbf{P}_{e}\mathbf{H}^{T} + \mathbf{R})^{-1}(\mathbf{D} - \mathbf{H}\mathbf{X}^{f}), \quad (3)$$

where **X** is an ensemble matrix of *n* model state realisations,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, ..., \mathbf{x}_n]$ ; and **D** is an ensemble matrix of observations. This ensemble is created by adding *n* realisations of random perturbations to the vector of observations **d**. The forecast error covariance matrix is approximated for the model of ensemble predictions by

$$\mathbf{P}_{e} = \frac{\mathbf{X}' \mathbf{X}'^{\mathrm{T}}}{n-1},\tag{3}$$

where  $\mathbf{X}'$  is a matrix of the ensemble perturbations of  $\mathbf{X}$  about a mean  $\overline{\mathbf{x}}$ . An ensemble approximation of observation error covariances is also possible, but has well understood problems (Kepert 2004), and is not employed here.

Because of sampling error introduced through the use of perturbed observations in  $\mathbf{D}$ , the EnKF formulation is expected to be less accurate than one that does not require perturbed observations (Whitaker and Hamill, 2002). In response to this Whitaker and Hamill (2002) proposed the Ensemble Square Root Filter, which is in most respects similar to the EnKF, but does not require perturbed observations.

After the analysis of observations the analysis error covariance should be

$$\mathbf{P}^{a} = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}^{f} (\mathbf{I} - \mathbf{K}\mathbf{H})^{\mathrm{T}} + \mathbf{K}\mathbf{R}\mathbf{K}^{\mathrm{T}}, \qquad (4)$$

where  $\mathbf{P}^{f}$  is the pre-analysis forecast error covariance. If the observations in equation (3) are not perturbed the post-analysis forecast error covariance  $\mathbf{P}^{a}$  will be underestimated, as the last term of equation (4) disappears (Whitaker and Hamill, 2002). While the EnSRF does not use perturbed observations it avoids underestimating the post-analysis forecast error covariance,  $\mathbf{P}^{a}$ , by updating the perturbation matrix  $\mathbf{X}'$  and the ensemble mean  $\bar{\mathbf{x}}$  separately. The ensemble mean is updated by the usual Kalman gain  $\mathbf{K}$  given in equation (2), while the perturbations are updated by a new gain matrix

$$\widetilde{\mathbf{K}} = \mathbf{P}^{f} \mathbf{H}^{\mathrm{T}} ((\sqrt{\mathbf{H}\mathbf{P}^{f}} \mathbf{H}^{\mathrm{T}} + \mathbf{R})^{-1})^{\mathrm{T}} \times (\sqrt{\mathbf{H}\mathbf{P}^{f}} \mathbf{H}^{\mathrm{T}} + \mathbf{R} + \sqrt{\mathbf{R}})^{-1} .$$
(5)

In this paper we consider the sequential assimilation application in both the one and three dimensional forms. In the one dimensional form, each observation is assimilated independently modifying the temperature of the water column directly beneath it. Only the ensemble perturbations from the model cells beneath the observation are used to generate the forecast error covariance matrix.

The advantage of the one dimensional approach is that by using single observations the **R** and  $\mathbf{HP}^{f}\mathbf{H}^{T}$ matrices collapse to scalars. This means that the computational cost of this assimilation form is significantly less than the three dimensional assimilation form. In contrast, the three dimensional form assimilates all the observations together to derive a correction to water temperature throughout the entire model domain. Solving this requires matrix inversions of large matrices (2000+ elements) but has the advantage that all the error covariance information is shared throughout the domain.

# 3 Experimental Setup

The data assimilation experiments described in this paper are illustrated through their application to a case study. The case study location is Port Phillip Bay, a shallow (~20m depth) highly enclosed basin located in south eastern Australia.

There are three sources of heating and cooling for the bay. These are the tidal exchange of water with the open sea through 'the heads', atmospheric heat fluxes, and riverine inputs of water. Of these the riverine influence is minimal and the atmospheric heat fluxes are the predominant heating and cooling mechanism.



Figure 3: Location of monitoring sites (circles) and weather stations (triangles) used in the modelling.

Tidal exchange is limited by the narrow entrance at 'the heads'.

The synthetic data assimilation experiments conducted here use a twin experiment technique to examine the improvement in prediction capability. The basis of this technique is to compare how close a degraded model approaches its 'true' twin when observations are assimilated.

The following set-up is used. An initial 'truth' model run was undertaken using atmospheric data collected predominantly at Point Wilson (see Figure 3). Data at Laverton were used for data types not collected at Point Wilson. Snapshots of the surface (top 1m) layer were extracted at an interval of two days to create a set of 'satellite' observations. These were degraded through the addition of spatially uncorrelated noise with a standard deviation of 0.3°C, which is in the range of the error associated with satellite observed SST (Brown and Minnett, 1999).

A second model run was set up using atmospheric forcing predominantly from Frankston (Figure 3) with data from Moorabbin for data types not collected at Frankston. Evaporation data was common to both model runs as only one station collected it. This different atmospheric forcing data was used to represent the typical uncertainty associated with model forcing data. The initial water temperature was set 1°C warmer to represent the uncertainty associated with specifying model initial conditions. This second model run, termed the 'open loop' run, indicates what prediction performance would be expected if there was no assimilation.

Four data assimilation runs were performed based on the second model set-up, with each run differing by the data assimilation technique and form used. The ensemble initialisation and propagation followed the procedure outlined in Turner et al (2005). In all cases 14 ensemble members were used as a trade-off between accuracy and computational time. In both the truth and open loop runs an initial 10 day model spinup was undertaken to produce realistic current fields from initially calm conditions, and the models run for a period of 40 days.

The hydrodynamic modelling was undertaken using the CSIRO Model for Estuaries and Coastal Oceans (MECO). MECO is a finite difference hydrodynamic model based on the three dimensional equations of momentum, continuity, and conservation of heat and salt, utilising the hydrostatic and Boussinesq assumptions (Walker et al, 2002).

The 'original' MECO thermodynamics formulation was used in all simulations. This formulation computes bulk values of the components of the energy balance. The net heat flux due to longwave and shortwave radiation together with sensible and latent heat flux is used to adjust the surface layer temperature in this formulation. Heating effects due to shortwave radiation absorption through the water column are also included.

An example of model output indicating currents superimposed over water elevation is presented in Figure 4. This figure shows some of the characteristics of Port Phillip Bay which initiated this study. The narrow entrance to the bay constricts the flow, as illustrated by the high velocities at the entrance. This limits exchange between Port Phillip Bay and Bass Strait to the southern portion of the bay; approximately the area in Figure 4 covered by the dense concentration of flow arrows.

While MECO allows for the use of spatially varying atmospheric inputs, spatially uniform inputs have been used throughout this paper. The justification for this is that the spatial extent is not so large as to warrant the complications of generating spatially varying wind fields, and other atmospheric inputs are not expected to vary significantly.

# 4 Results

To illustrate the experiment results, time series of water temperature have been extracted from three comparison sites (see Figure 3) at different locations within the water column in each case (Figure 4). These figures demonstrate that the data assimilation is able to positively impact the deeper unobserved layers of the water column in addition to the observed surface layer. Each figure compares the open loop and ensemble mean predicted by the different data assimilation techniques with the truth.

Consider first the difference between the truth and the open loop in each of the three cases.

The initial difference in water temperature prediction is due to the 1°C rise made to the degraded model initial condition. Over time the graphs mirror each other and. move closer. This is an effect of the initial condition error being diminished as a result of little or no error in the open boundary and atmospheric forcing, and the same model physics (heating and cooling) used for both scenarios. While atmospheric



Figure 4: Sample of output from MECO model. Shading indicates water elevation, while arrows indicate surface currents during an ebb tide.

data sets are taken from different locations they are not so different to produce wildly differing predictions.

For all three comparison locations the introduction of observations through data assimilation significantly improves the model water temperature forecast. Any initial discrepancy between the open loop and the ensemble mean is due to the ensemble generation technique. The initial assimilation generates significant improvements in prediction with subsequent assimilation times resulting in smaller corrections. As the observation error is constant throughout the assimilation period, the amount of improvement made by the assimilation is dependent upon the uncertainty associated with the model prediction. This in turn is a function of the spread of the ensemble members about the ensemble mean and is calculated using equation (3).

A comparison between the different filtering techniques and open loop run is made by contrasting the root mean squared (RMS) difference between the ensemble mean forecast and the truth for the entire model domain (Figure 5). This shows that the initial average forecast error is about 0.8°C; significantly worse than the prescribed observation error of 0.3°C; the initial 1°C initial condition difference was reduced during the spin-up period. Moreover, assimilation significantly improved the model predictions both in terms of absolute RMS and relative to the open loop.

Contrasting the one dimensional examples; the EnSRF performed significantly better than the EnKF. The poor performance of the one dimensional EnKF is due to sampling error in the perturbed observations, which is magnified by the one dimensional form. Using a three dimensional assimilation form gave significantly better results than the one dimensional assimilation form, a consequence of more information being available to the three dimensional form. Both techniques appeared equally as effective in the three dimensional form, although the EnSRF performed slightly better over the first 10 days of the assimilation period. Sampling error in the EnKF appears to be reduced by averaging over time.

The improvements in prediction can be calculated relative to the open loop prediction (Figure 5) by



Figure 4: Time series of predicted water temperature for the open loop and assimilated model runs as compared to the truth run for (a) Long Reef monitoring location for the surface (depth 1m); (b) Hobsons Bay monitoring location for the middle (depth 5m); and (c) Central monitoring location for the bottom (depth 23m).

dividing the difference between the RMS error of the assimilation and the RMS error of the open loop by the RMS error of the open loop prediction. Both three dimensional assimilation techniques reached 90% improvement over the assimilation period, while the one dimensional EnSRF averaged 80-90% and the one



Figure 5: RMS difference between ensemble mean prediction and truth for the entire model domain.

dimensional EnKF averaged 60-70% over the assimilation period.

#### 5 Discussions

The results obtained through the twin experiments show that the three dimensional form performed significantly better than the one dimensional form. This result is not unexpected. By using a three dimensional form the filter is able use the spatial relationships inherent in the model forecast to produce an improved update.

Although both filter types performed well in the three dimensional form, the EnSRF performed slightly better. It had a better prediction for the first 10 days of the assimilation period and as good prediction for the remainder of the period. Moreover, it predicted the estimate of RMS error slightly better than the EnKF. This is due to the small number of ensembles (14 in this application), which in combination with the sampling error inherent in the EnKF reduces its performance. With a longer assimilation run or more ensemble members the EnKF and EnSRF are expected to give equivalent results, which was shown with the three dimensional EnKF approaching the EnSRF over time (Figure 5). However, the number of ensemble members is the variable which most significantly influences computational costs, with the assimilation step being relatively inexpensive. Still, these tests cannot be considered thorough enough to state with certainty that the EnSRF is the better choice.

The experiments have shown that significant improvements to model predictions of water temperature are possible through the assimilation of surface layer observations. These experiments however, have been performed in an artificial environment and the same level of performance cannot be expected in a real case. One of the main limitations of this study is that the same model equations were used to create the observations and truth data as for the assimilation and open loop forecasts. This means the model was conditioned to perform well relative to the truth. In reality, the physical process occurring cannot be modelled exactly and in a real case more errors will be introduced in this way.

### 6 Conclusions

Two sequential data assimilation techniques have been compared in both one and three dimensional forms. While the three dimensional form is superior to the one dimensional form, the EnSRF performs only slightly better than the EnKF and the EnSRF is less affected by dimensionality. While the performance of both filters is admirable, an application using real observations is necessary to better appreciate the improvements actually realised from data assimilation.

This paper has demonstrated the significant potential of data assimilation techniques to improve the reliability and accuracy of a model prediction. All of the assimilation techniques used gave significant improvement over the control without assimilation. The techniques described are, more generally, applicable to other coastal marine modelling settings and will improve any coastal forecasting system.

### 7 Acknowledgements

The work undertaken in this paper has been funded by a University of Melbourne CSIRO Collaborative Grant and a Melbourne Research Scholarship. The modelling data was provided by Melbourne Water, Melbourne Ports Authority, Bureau of Meteorology and MAFRI.

#### 8 References

Brown, O.B. and P.J. Minnett (1999) MODIS infrared sea surface temperature algorithm theoretical basis document, Version 2, University of Miami.

Evensen, G. (2003) The ensemble Kalman filter: theoretical formulation and practical implementation, *Ocean Dynamics*, Vol 53, 343-367.

Ghil, M. and P. Malanotte-Rizoli (1991) Data assimilation in meteorology and oceanography, *Advances in Geophysics*, Vol 33, 141-266.

Houtekamer, P.L and H.L. Mitchell (1998) Data assimilation using an ensemble Kalman filter technique, *Monthly Weather Review*, Vol 126, 796-811.

Kepert, J.D. (2004) On ensemble representation of the observation-error covariance in the ensemble Kalman filter, Ocean Dynamics, (In Press)

Turner, M.R.J., J.P. Walker and P.R. Oke, (In preparation) Ensemble member generation for sequential data assimilation, *Ocean Modelling*.

Walker, S.J., J.R. Waring, M. Herzfeld and P. Sakov (2002) *Model for Estuaries and Coastal Oceans: User Manual*, Version 4.01, CSIRO Marine, Hobart.

Whitaker, J.S and T.H. Hamill (2002) Ensemble data assimilation without perturbations. *Monthly Weather Review*, Vol. 130, 1913–1924.