# Model-based Prognostics with Fixed-lag Particle Filters

**Matthew Daigle** [1] **and Kai Goebel** [2]

[1] *University of California, Santa Cruz, NASA Ames Research Center, Moffett Field, CA, 94035, USA*
*matthew.j.daigle@nasa.gov*
[2] *NASA Ames Research Center, Moffett Field, CA, 94035, USA*
*kai.goebel@nasa.gov*

## ABSTRACT

Model-based prognostics exploits domain knowledge of the system, its components, and how they fail by casting the underlying physical phenomena in a physics-based model that is derived from first principles. In most applications, uncertainties from a number of sources cause the predictions to be inaccurate and imprecise even with accurate models. Therefore, algorithms are employed that help in managing these uncertainties. Particle filters have become a popular choice to solve this problem due to their wide applicability and ease of implementation. We present a general model-based prognostics methodology using particle filters. In order to provide more accurate and precise estimates, and, therefore, more accurate and precise predictions, we investigate the use of fixed-lag filters. We develop a detailed physics-based model of a pneumatic valve, and perform comprehensive simulation experiments to illustrate our prognostics approach. The experiments demonstrate the advantages that fixed-lag filters may provide in the context of prognostics, as measured by prognostics performance metrics.

## 1 INTRODUCTION

Prognostics is a key enabling technology for applying condition-based maintenance. The goal of prognostics is to make *end of life* (EOL) and *remaining useful life* (RUL) predictions that enable timely maintenance decisions to be made. As with diagnostics, prognostics methods may typically be categorized as either data-driven or model-based approaches. Data-driven prognostics approaches rely on run-to-failure data that are used to train algorithms to recognize trends and estimate EOL and RUL. Indeed, data-driven prognostic approaches dominate the literature at this point; see (Schwabacher, 2005) for a survey. However, there are numerous cases where the necessarily large amount of run-to-failure data does not exist.

Here, model-based approaches offer a viable alternative. Model-based prognostics approaches exploit domain knowledge of the system, its components, and how they fail in order to provide accurate EOL and RUL predictions (Roemer *et al.*, 2005; Byington *et al.*, 2004; Saha and Goebel, 2009). The underlying physical phenomena are captured in a physics-based model that is derived from first principles, therefore, model-based approaches can provide EOL and RUL estimates that are much more accurate and precise than data-driven approaches, if the models are accurate. Still, modeling the physics of a system (or even just a component) is rarely a trivial task.

We adopt a model-based prognostics approach that is based on joint state-parameter estimation. Many model-based prognostics frameworks perform state and parameter estimation using particle filters, which approximate the posterior as a set of discrete, weighted samples. Although suboptimal, the advantage of particle filters is that they can be applied to systems which may be nonlinear and have non-Gaussian noise terms, where optimal solutions are unavailable or intractable. Further, because they are based on probability distributions, they help in managing the uncertainty that may arise from a number of sources. In (Saha and Goebel, 2009), the authors apply a particle filtering approach to prediction of end of discharge and EOL in lithium-ion batteries. In (Orchard *et al.*, 2008), the authors present a particle filter-based diagnosis and prognosis framework using correction loops, with application to crack growth in aircraft components. In (Abbas *et al.*, 2007), the authors apply a particle filter-based prognosis method to prediction of battery grid corrosion.

Similar to these approaches, we also develop a general model-based prognostics methodology using particle filters. Unlike previous work, however, we investigate the use of *fixed-lag* filters to improve estimation and, subsequently, prediction. Fixed-lag filters incorporate observations beyond a given time point to calculate the estimate at that time point. Since more information is being used, estimates can be more accurate and more precise. The disadvantage to using fixed-lag filters is that state estimates are delayed, i.e., the state estimate for time $t$ is only computed once observations from some later time are available. Since prognostics often deals with very large time scales, the delay in the state estimate inherent with fixed-lag filters is acceptable given the possible improvements in
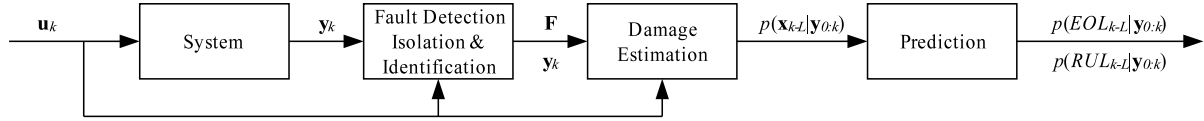
Figure 1: Prognostics architecture.

estimation and prediction that they offer. However, to our knowledge, the use of fixed-lag filters in the context of prognostics has not been explored.

In this paper, we develop a prognostics framework that incorporates fixed-lag particle filters. As a case study, we construct a detailed physics-based model of a pneumatic valve, and use this model to study the effects of different damage mechanisms. We run a number of prognostics experiments in simulation to demonstrate how fixed-lag filters may improve estimation and EOL/RUL predictions. Prognostics performance is evaluated using established prognostic metrics (Saxena *et al.*, 2008; 2009).

The paper is organized as follows. Section 2 formulates the prognostics problem and overviews the computational architecture we adopt. Section 3 develops the damage estimation method using fixed-lag particle filters. Section 4 describes the EOL/RUL prediction procedure. Section 5 presents the pneumatic valve case study with experimental results in simulation. Section 6 concludes the paper.

## 2 PROGNOSTICS APPROACH

### 2.1 Problem Formulation

The problem of prognostics is to predict the EOL and/or the RUL of a component, where EOL is defined as the time point at which a component no longer meets specified functional and/or performance requirements, and RUL is the time remaining until that point. In this paper, we develop a general model-based approach, where the system model is given by

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t), \mathbf{v}(t))$$
$$\mathbf{y}(t) = \mathbf{h}(t, \mathbf{x}(t), \mathbf{u}(t), \mathbf{n}(t)),$$

where $\mathbf{x}(t) \in \mathbb{R}^{n_\mathbf{x}}$ is the state vector, $\mathbf{u}(t) \in \mathbb{R}^{n_\mathbf{u}}$ is the input vector, $\mathbf{v}(t) \in \mathbb{R}^{n_\mathbf{v}}$ is the process noise vector, $\mathbf{f}$ is the state equation, $\mathbf{y}(t) \in \mathbb{R}^{n_\mathbf{y}}$ is the output vector, $\mathbf{n}(t) \in \mathbb{R}^{n_\mathbf{n}}$ is the measurement noise vector, and $\mathbf{h}$ is the output equation.

Our goal is to predict EOL at a given time point $t_P$ using the discrete sequence of observations up to time $t_P$, denoted as $\mathbf{y}_{0:t_P}$. In order to determine when EOL has been reached, we require a condition that is a function of the system state, $C_{EOL}(\mathbf{x}(t))$, which determines whether EOL has been reached, where

$$C_{EOL}(\mathbf{x}(t)) = \begin{cases} 1, & \text{if EOL is reached} \\ 0, & \text{otherwise.} \end{cases}$$

Using this function, we can define EOL as

$$EOL(t_P) \triangleq \underset{t \geq t_P}{\arg\min} \, C_{EOL}(\mathbf{x}(t)) = 1.$$

RUL is then simply defined as

$$RUL(t_P) \triangleq EOL(\mathbf{x}(t_P)) - t_P.$$

Because of the noise inherent in the process and the measurements, we must compute a probability distribution of the EOL or RUL, i.e., the goal is to compute, at time $t_P$, $p(EOL(t_P)|\mathbf{y}_{0:t_P})$ or $p(RUL(t_P)|\mathbf{y}_{0:t_P})$.

### 2.2 Prognostics Architecture

We adopt a model-based approach, wherein we develop detailed physics-based models of components and systems that include descriptions of how fault parameters evolve in time. These models depend on unknown and possibly time-varying damage/wear parameters $\boldsymbol{\theta}(t) \subset \mathbf{x}(t)$. Therefore, our solution to the prognostics problem takes the perspective of joint state-parameter estimation. In discrete time $k$, we estimate $\mathbf{x}_k$ and use the estimates to predict EOL and RUL at desired time points.

In order to improve the state estimates, and, therefore predictions, we utilize fixed-lag filters, where we compute $p(\mathbf{x}_{k-L}|\mathbf{y}_{0:k})$, where $L$ is the lag. Using $p(\mathbf{x}_{k_P-L}|\mathbf{y}_{0:k_P})$ at time $k_P - L$, we compute $p(EOL_{k_P-L}|\mathbf{y}_{0:k_P})$ and $p(RUL_{k_P-L}|\mathbf{y}_{0:k_P})$.

We employ the prognostics architecture in Fig. 1. The system is provided with inputs $\mathbf{u}_k$ and provides measured outputs $\mathbf{y}_k$. The fault detection, isolation, and identification (FDII) module provides a fault set $\mathbf{F}$, which is used by the damage estimation module to determine estimates of the states including unknown parameters, represented as a probability distribution $p(\mathbf{x}_{k-L}|\mathbf{y}_{0:k})$. This distribution is used by the prediction module, which computes EOL and RUL using hypothesized future inputs. EOL and RUL are computed as probability distributions $p(EOL_{k_P-L}|\mathbf{y}_{0:k_P})$ and $p(RUL_{k_P-L}|\mathbf{y}_{0:k_P})$. In this paper, we focus on the damage estimation and prediction modules, and assume a solution to FDII.

## 3 DAMAGE ESTIMATION

To estimate the damage, we need to estimate $p(\mathbf{x}_k|\mathbf{y}_{0:k})$. In this paper, we use the particle filter for this purpose (Arulampalam *et al.*, 2002; Cappe *et al.*, 2007). With particle filters, the state distribution is approximated by a set of discrete weighted samples, or particles, $\{\mathbf{x}_k^i, w_k^i\}_{i=1}^N$, where $N$ denotes the number of particles, $\mathbf{x}_k^i$ denotes the state estimate for particle $i$, and $w_k^i$ denotes the weight of particle $i$.

Particle filters are best suited to estimation in nonlinear systems with possibly non-Gaussian noise, where optimal solutions are unavailable or intractable. In this respect, they can be viewed as a general (suboptimal) solution to the state estimation problem. Performance can be improved by increasing the number of particles, but this also results in higher computational costs. The number of particles must be chosen to suit the application requirements.

As described in Section 2, parameters augment the state vector, i.e., $\boldsymbol{\theta}_k \subset \mathbf{x}_k$. In this way, the particle filter is being used to perform joint state-parameter estimation. Here, the parameters $\boldsymbol{\theta}_k$ evolve by some unknown random process that is independent of the state $\mathbf{x}_k$. To perform parameter estimation within a particle filter framework, however, we need to assign some type of evolution to the parameters. The typical solution is to use a random walk, i.e., for parameter $\theta$, $\theta_k = \theta_{k-1} + \xi_{k-1}$, where $\xi_{k-1}$ is a noise term, and typically Gaussian. During the sampling step, particles are generated with parameter values that will be different from the initial guesses for the unknown parameters. The particles with parameter values closest to the true values should be assigned higher weight, thus allowing the particle filter to converge to the true values. The selected variance of the random walk noise must be large enough so as to allow convergence in a reasonable amount of time, but small enough such that when convergence is reached, the parameter can be tracked smoothly. Since the parameter values are unknown to start with, this can be a difficult task, but knowledge of the correct order of magnitude of the parameter is helpful. If the unknown parameters are constant, then other approaches can be employed to improve estimates and offset the increase in covariance contributed by the random walk (Liu and West, 2001; Clapp and Godsill, 1999).

We employ the *sampling importance resampling* (SIR) particle filter, and implement the resampling step using systematic resampling (Kitagawa, 1996). In particle filters, the posterior density is approximated by

$$p(\mathbf{x}_k|\mathbf{y}_{0:k}) \approx \sum_{i=1}^{N} w_k^i \delta(\mathbf{x}_k - \mathbf{x}_k^i).$$

The *fixed-lag smoothing* distribution can be approximated by (Doucet *et al.*, 2000)

$$p(\mathbf{x}_{k-L}|\mathbf{y}_{0:k}) \approx \sum_{i=1}^{N} w_k^i \delta(\mathbf{x}_{k-L} - \mathbf{x}_{k-L}^i).$$

So, to compute $p(\mathbf{x}_{k-L}|\mathbf{y}_{0:k})$, we need to determine the weight of the particle at time $k$. This can be computed by running the standard particle filter algorithm within an inner loop up to time $k$ to determine the weight of the particles at time $k$, and assigning these weights to the particles for time $k - L$. The particle filter may then proceed from time $k - L$ as usual. However, the resampling step complicates this procedure. During resampling, particles may be either dropped or multiplied. This results in a loss of diversity of the particle paths and smoothed estimates based on these paths degnerate (Arulampalam *et al.*, 2002; Cappe *et al.*, 2007). The solution we adopt is to skip the resampling step during the lookahead portion, which avoids the degeneracy that would be introduced by resampling.

The pseudocode for a single step of the fixed-lag SIR filter is shown as Algorithm 1. Each particle is propagated forward to time $k$ (without resampling), and the particle weight is assigned using $y_k$. As with the standard SIR filter, the weights are then normalized, followed by the resampling step (see (Arulam-

---

**Algorithm 1** Fixed-lag SIR Filter

**Inputs:** $\{\mathbf{x}_{k-L-1}^i, w_{k-L-1}^i\}_{i=1}^N, \mathbf{u}_{k-L:k}, \mathbf{y}_k$
**Outputs:** $\{\mathbf{x}_{k-L}^i, w_{k-L}^i\}_{i=1}^N$
**for** $i = 1$ **to** $N$ **do**
  **for** $j = 0$ **to** $L$ **do**
    $\mathbf{x}_{k-L+j}^i \sim p(\mathbf{x}_{k-L+j}|\mathbf{x}_{k-L+j-1}^i, \mathbf{u}_{k-L+j-1})$
  **end for**
  $w_{k-L}^i \leftarrow p(\mathbf{y}_k|\mathbf{x}_k^i, \mathbf{u}_k)$
**end for**
$W \leftarrow \sum_{i=1}^N w_{k-L}^i$
**for** $i = 1$ **to** $N$ **do**
  $w_{k-L}^i \leftarrow w_{k-L}^i/W$
**end for**
$\{\mathbf{x}_{k-L}^i, w_{k-L}^i\}_{i=1}^N \leftarrow \texttt{Resample}(\{\mathbf{x}_{k-L}^i, w_{k-L}^i\}_{i=1}^N)$

---

palam *et al.*, 2002) for pseudocode). With $L = 0$, the algorithm is equivalent to the standard SIR filter.

## 4 PREDICTION

In the prediction phase, we wish to compute at time $k_P - L$, the distributions $p(EOL_{k_P-L}|\mathbf{y}_{0:k_P})$ and $p(RUL_{k_P-L}|\mathbf{y}_{0:k_P})$. The fixed-lag particle filter computes

$$p(\mathbf{x}_{k_P-L}|\mathbf{y}_{0:k_P}) \approx \sum_{i=1}^{N} w_{k_P-L}^i \delta(\mathbf{x}_{k_P-L} - \mathbf{x}_{k_P-L}^i).$$

We can approximate a prediction distribution $n$ steps forward as (Doucet *et al.*, 2000)

$$p(\mathbf{x}_{k_P-L+n}|\mathbf{y}_{0:k_P}) \approx$$
$$\sum_{i=1}^{N} w_{k_P-L}^i \delta(\mathbf{x}_{k_P-L+n} - \mathbf{x}_{k_P-L+n}^i).$$

So, for a given state $\mathbf{x}_{k_P-L}^i$ propagated $n$ steps forward (without new data), we can simply take its weight as $w_{k_P-L}^i$. Similarly, we can approximate the EOL as

$$p(EOL_{k_P-L}|\mathbf{y}_{0:k_P}) \approx$$
$$\sum_{i=1}^{N} w_{k_P-L}^i \delta(EOL_{k_P-L} - EOL_{k_P-L}^i).$$

The idea, then, is to propagate each particle forward to EOL and use the particle's weight at time $k_P - L$ for the weight of the EOL prediction.

The pseudocode for the prediction procedure is given as Algorithm 2. Each particle $i$ is propagated forward until $C_{EOL}(\mathbf{x}_k^i)$ evaluates to 1, at which EOL has been reached for this particle. Prediction requires hypothesizing future inputs of the system $\hat{\mathbf{u}}_k$. The inputs must be chosen carefully because different inputs often have different effects on damage progression. The choice depends on the particular application.

## 5 CASE STUDY

In order to illustrate our prognostics methodology, we take a pneumatic valve as a case study. We develop

**Algorithm 2** EOL Prediction

---

**Inputs:** $\{\mathbf{x}^i_{k_P-L}, w^i_{k_P-L}\}^N_{i=1}$
**Outputs:** $\{EOL^i_{k_P-L}, w^i_{k_P-L}\}^N_{i=1}$
**for** $i = 1$ **to** $N$ **do**
    $k \leftarrow k_P - L$
    $\mathbf{x}^i_k \leftarrow \mathbf{x}^i_{k_P-L}$
    **while** $C_{EOL}(\mathbf{x}^i_k) = 0$ **do**
        Predict $\hat{\mathbf{u}}_k$
        $\mathbf{x}^i_{k+1} \sim p(\mathbf{x}_{k+1}|\mathbf{x}^i_k, \hat{\mathbf{u}}_k)$
        $\mathbf{x}^i_k \leftarrow \mathbf{x}^i_{k+1}$
        $k \leftarrow k + 1$
    **end while**
    $EOL^i_{k_P-L} \leftarrow k$
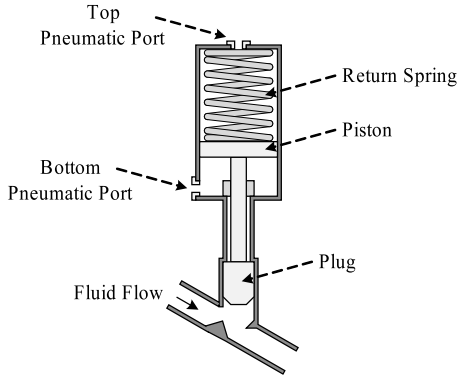**end for**

---



Figure 2: Pneumatic valve.

a physics-based model of the valve and its damage mechanisms. We then present simulation experiments to demonstrate parameter estimation and EOL and RUL prediction for different choices of the lag $L$.

## 5.1 Component Modeling

Pneumatic valves are complex mechanical systems used in many domains. These valves are actuated by gas, and can use different types of actuators. A normally-closed valve with a linear cylinder actuator is depicted in Fig. 2. The valve is opened by filling the chamber below the piston with gas up to the supply pressure, and evacuating the chamber above the piston down to atmospheric pressure. The valve is closed by filling the chamber above the piston, and evacuating the chamber below the piston. The return spring ensures that when pressure is lost, the valve will close due to the force exerted by the return spring.

We develop a physics model of the valve based on mass and energy balances. The system state includes the position of the valve, $x(t)$, the velocity of the valve, $v(t)$, the mass of the gas in the volume above the piston, $m_t(t)$, and the mass of the gas in the volume below the piston, $m_b(t)$:

$$\mathbf{x}(t) = \begin{bmatrix} x(t) \\ v(t) \\ m_t(t) \\ m_b(t) \end{bmatrix}.$$

The position when the valve is fully closed is defined as $x = 0$. The stroke length of the valve is denoted by $L_s$; when the valve is fully open its position is $x = L_s$.

The derivatives of the states are described by

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} v(t) \\ \frac{1}{m}\sum F \\ f_t(t) \\ f_b(t) \end{bmatrix},$$

where $m$ is the combined mass of the piston and plug, $\sum F$ is the sum of forces acting on the valve, and $f_t(t)$ and $f_b(t)$ are the mass flows going into the top and bottom pneumatic ports, respectively.

The inputs are considered to be

$$\mathbf{u}(t) = \begin{bmatrix} p_l(t) \\ p_r(t) \\ u_t(t) \\ u_b(t) \end{bmatrix},$$

where $p_l(t)$ and $p_r(t)$ are the fluid pressures on the left and right side of the plug, respectively, and $u_t(t)$ and $u_b(t)$ are the input pressures to the top and bottom pneumatic ports. These pressures will alternate between the supply pressure and atmospheric pressure depending on the commanded valve position.

The sum of forces acting on the piston includes (1) the forces from the pneumatic gas: $(p_b(t) - p_t(t))A_p$, where $p_b(t)$ and $p_t(t)$ are the gas pressures on the bottom and the top, respectively, and $A_p$ is the surface area of the piston, (2) the forces from the fluid flowing through the valve: $(p_r(t) - p_l(t))A_v$, where $A_v$ is the area of the valve contacting the fluid, (3) the weight of the moving parts of the valve: $-mg$, where $g$ is the acceleration due to gravity, (4) the spring force: $-k(x(t) - x_o)$, where $k$ is the spring constant and $x_o$ is the amount of spring compression when the valve is closed, (5) friction: $-rv(t)$, where $r$ is the coefficient of kinetic friction, and (6) the contact forces at the boundaries of the valve motion:

$$\begin{cases} k_c(-x), & x < 0 \\ 0, & 0 \leq x \leq L_s \\ -k_c(x - L_s), & x > L_s, \end{cases}$$

where $k_c$ is the (large) spring constant associated with the flexible seals.

The pressures $p_t(t)$ and $p_b(t)$ are calculated as:

$$p_t(t) = \frac{m_t(t)R_gT}{V_{t_0} + A_p(L_s - x(t))}$$

$$p_b(t) = \frac{m_b(t)R_gT}{V_{b_0} + A_px(t)}$$

where we assume an isothermal process in which the gas temperature is constant at $T$, $R_g$ is the gas constant for the pneumatic gas, and $V_{t_0}$ and $V_{b_0}$ are the minimum gas volumes for the gas chambers above and below the piston, respectively.

The gas flows are given by:

$$f_t(t) = f_g(p_t(t), u_t(t))$$
$$f_b(t) = f_g(p_b(t), u_b(t))$$

where $f_g$ defines gas flow through an orifice for choked and non-choked flow conditions (Perry and Green, 2007): $f_g(p_1, p_2) =$

$$
\begin{cases}
C_s A_s p_1 \sqrt{\frac{\gamma}{Z R_g T} \left(\frac{2}{\gamma+1}\right)^{\frac{\gamma+1}{\gamma-1}}}, \\
\qquad p_1 \geq p_2 \wedge p_1/p_2 \geq \left(\frac{\gamma+1}{2}\right)^{\gamma/(\gamma-1)} \\
C_s A_s p_1 \sqrt{\frac{\gamma}{Z R_g T} \left(\frac{2}{\gamma-1}\right) \left(\left(\frac{p_2}{p_1}\right)^{\frac{2}{\gamma}} - \left(\frac{p_2}{p_1}\right)^{\frac{\gamma+1}{\gamma}}\right)}, \\
\qquad p_1 \geq p_2 \wedge p_1/p_2 < \left(\frac{\gamma+1}{2}\right)^{\gamma/(\gamma-1)} \\
C_s A_s p_2 \sqrt{\frac{\gamma}{Z R_g T} \left(\frac{2}{\gamma+1}\right)^{\frac{\gamma+1}{\gamma-1}}}, \\
\qquad p_1 < p_2 \wedge p_2/p_1 \geq \left(\frac{\gamma+1}{2}\right)^{\gamma/(\gamma-1)} \\
C_s A_s p_2 \sqrt{\frac{\gamma}{Z R_g T} \left(\frac{2}{\gamma-1}\right) \left(\left(\frac{p_1}{p_2}\right)^{\frac{2}{\gamma}} - \left(\frac{p_1}{p_2}\right)^{\frac{\gamma+1}{\gamma}}\right)}, \\
\qquad p_1 < p_2 \wedge p_2/p_1 < \left(\frac{\gamma+1}{2}\right)^{\gamma/(\gamma-1)}
\end{cases}
$$

where $\gamma$ is the ratio of specific heats, $Z$ is the gas compressibility factor, $C_s$ is the flow coefficient, and $A_s$ is the orifice area. Choked flow occurs when the pressure ratio exceeds $\left(\frac{\gamma+1}{2}\right)^{\gamma/(\gamma-1)}$.

We select our measurement vector as

$$
\mathbf{y}(t) = \begin{bmatrix} x(t) \\ p_t(t) \\ p_b(t) \\ f_v(t) \end{bmatrix}
$$

where $f_v$ is the fluid flow through the valve:

$$
f_v(t) = \frac{x(t)}{L_s} C_v A_v \sqrt{\frac{2}{\rho} |p_{fl} - p_{fr}|} \, sign(p_{fl} - p_{fr}),
$$

where $C_v$ is the (dimensionless) flow coefficient of the valve, and $\rho$ is the liquid density, and we assume a linear flow characteristic for the valve.

Fig. 3 shows a nominal valve cycle. The valve is commanded to open at $0$ s. The top pneumatic port opens to atmosphere and the bottom opens to the supply pressure (approximately $5.3$ MPa, or $750$ psig). When the force on the underside of the piston is large enough to overcome the return spring, friction, and the gas force on the top of the piston, the valve begins to move upward as the pneumatic gas continues to flow into and out of the valve actuator. At about $8$ s the valve is completely open. The valve is commanded to close at $15$ s. The bottom pneumatic port opens to atmosphere and the bottom opens to the supply pressure. When the force balance becomes negative, the valve starts to move downward, and completely closes at around $20$ s. The valve closes faster than it opens due to the return spring.

## 5.2 Damage Modeling

Our general approach to damage modeling is as follows. First, we identify parameters in the model that characterize the extent of specific forms of damage, and these augment the state vector $\mathbf{x}$. We then incorporate models of how those parameters change over time with system operation. It is the parameters of
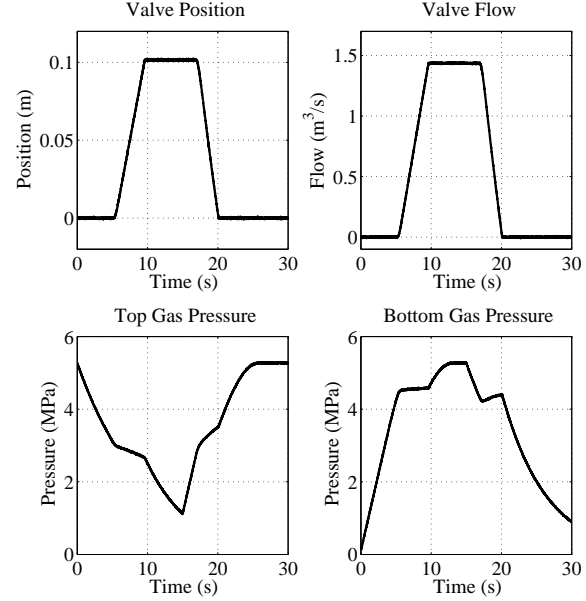


Figure 3: Nominal valve operation.

these equations that are unknown and must be estimated, and also augment $\mathbf{x}$ for that purpose. In the valve model, we consider damage or wear characterized by the increase in friction coefficient, the decrease in spring constant, the appearance and growth of an internal valve leak between the volumes on either side of the piston, and the appearance and growth of external leaks at the pneumatic ports.

One damage mechanism present in valves is sliding wear. The equation for sliding wear takes on the following form (Hutchings, 1992):

$$
\dot{V}(t) = w |F(t) v(t)|,
$$

where $V(t)$ is the wear volume, $w$ is the wear coefficient (which depends on material properties such as hardness), $F(t)$ is the sliding force, and $v(t)$ is the sliding velocity. Friction will increase linearly with sliding wear, because the contact area between the sliding bodies becomes greater as surface asperities wear down (Hutchings, 1992). Lubrication between the sliding bodies can also degrade over time. We therefore model the change in friction coefficient in a form similar to sliding wear:

$$
\dot{r}(t) = w_r |F_f(t) v(t)|
$$

where $w_r$ is the wear coefficient, and $F_f(t)$ is the friction force defined in the previous subsection. Fig. 4 shows the effect of an increase in friction on the valve cycle. From the simulation, we can determine the value of the friction parameter, $r^*$, at which the valve has reached EOL. At this value, the friction force becomes large enough that the valve cannot open within the $15$ s limit, as shown in Fig. 4. So, $C_{EOL}(\mathbf{x}(t)) = 1$ if $r(t) \geq r^*$.

We assume the same equation form for spring damage:

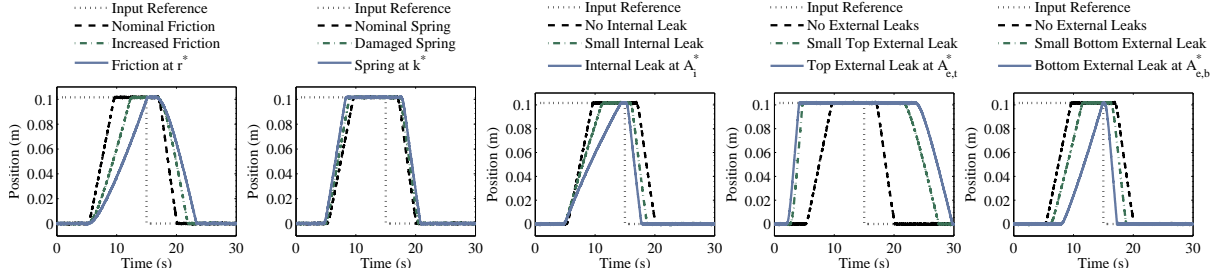$$
\dot{k}(t) = -w_k |F_s(t) v(t)|,
$$

Figure 4: Valve operation with increasing wear.

where $w_k$ is the spring wear coefficient and $F_s(t)$ is the spring force. The more the spring is used, the weaker it becomes. We define $k^*$ as the largest value of $k$ at which the valve will not fully close upon loss of supply pressure. Fig. 4 shows the effect of a decrease in the spring parameter on the valve cycle. In normal operation, without the spring tending the valve to close, the valve will open faster and close slower. However, the spring must be strong enough to close the valve against system pressure when the actuating pressure is lost. So, $C_{EOL}(\mathbf{x}(t)) = 1$ also if $k(t) \leq k^*$.

An internal leak in the valve can appear at the seal surrounding the piston as a result of sliding wear. The pneumatic gas is then able to flow between the volumes above and below the piston, decreasing the response time of the valve. We parameterize this leak by its equivalent orifice area, $A_i(t)$, described by:

$$\dot{A}_i(t) = w_i |F_f(t)v(t)|,$$

where $w_i$ is the wear coefficient. The mass flow at the leak, $f_i(t)$, is computed using the gas flow equation:

$$f_i(t) = f_g(p_t(t), p_b(t)).$$

As sliding wear occurs, the leak size keeps increasing. The presence of an internal leak makes it more difficult to actuate the valve, because it causes gas to flow into the lower pressure volume that is being evacuated and out of the higher pressure volume that is being filled. We define $A_i^*$ as the minimum internal leak area at which the valve cannot open within the 15 s limit. So, $C_{EOL}(\mathbf{x}(t)) = 1$ also if $A_i(t) \geq A_i^*$. Fig. 4 shows the effect of an internal leak on the valve cycle.

External leaks can also form, most likely at the actuator connections to the pneumatic gas supply, due to corrosion and other environmental factors. Without knowledge of how the leak size progresses, we assume the growth of the area of the leak holes, $A_e(t)$, is linear:

$$\dot{A}_e(t) = w_e,$$

where $w_e$ is the wear coefficient. We use additional $t$ and $b$ subscripts to denote leaks at the top and bottom pneumatic ports, respectively. The effect of the formation of a leak at the top pneumatic port is that it becomes easier to open the valve but more difficult to close it. Conversely, the effect of a leak at the bottom pneumatic port is that it becomes more difficult to open but easier to close the valve. Through simulation we can determine the minimum size leak holes at which the valve cannot open or close within

the 15 s limit, $A_{e,t}^*$ and $A_{e,b}^*$. (An alternative is to use a maximum allowable leakage rate to define EOL.) So, $C_{EOL}(\mathbf{x}(t)) = 1$ also if $A_{e,t}(t) \geq A_{e,t}^*$ or $A_{e,b}(t) \geq A_{e,b}^*$. Fig. 4 shows the effect of external leak on the valve cycle.

### 5.3 Experimental Results

We performed a number of simulation experiments to validate our prognostics methodology and evaluate the usefulness of fixed-lag filters for prognostics using performance metrics described in (Saxena *et al.*, 2008; 2009). In each experiment, we considered additive zero-mean process and measurement noise, used $N = 500$ particles, and used a sample time of 0.01 s. We assumed that only a single damage mechanism was active, and, in each experiment, started from the point where damage has been identified and the only unknown is the wear coefficient (initially assumed to be 0 for parameter estimation). We tuned our particle filters by adjusting the amount of process and measurement noise it considered, assuming the order of magnitude of the wear coefficients were known. With the noise variances selected, we then varied only the lag $L$.

**Validation of the Methodology**
First, we provide an example scenario that demonstrates the effectiveness of our model-based methodology. Fig. 5 shows estimation results for the case of an internal leak with $L = 3$. There was little error present in the tracking of the outputs, so we show only the estimation results of the internal leak area, $A_i(t)$, and its wear coefficient $w_i$. In our framework, accurate and precise tracking of the hidden damage parameter translates to accurate and precise predictions of EOL and RUL. In this case, the estimate of $w_i$ converges in about 3 cycles, or 90 s. After that point, $A_i(t)$ can be tracked well.

The EOL predictions for each prediction point (every 10 cycles, where one cycle corresponds to 30 s) are shown in Fig. 6, with a mixture of Gaussian distributions fitted to the particle populations for visualization purposes. The true EOL for the chosen value of $w_i$ is 106 cycles. The probability distributions all cover the true EOL, and as time progresses, the predictions become significantly more accurate and precise.

This result is also shown by the $\alpha$-$\lambda$ accuracy metric (Saxena *et al.*, 2008), as given in Fig. 7. Here, $\alpha \in [0, 1]$ defines bounds as a function of RUL, and $\lambda \in [0, 1]$ denotes the fraction of the time from the first prediction to the true EOL. We use the extended
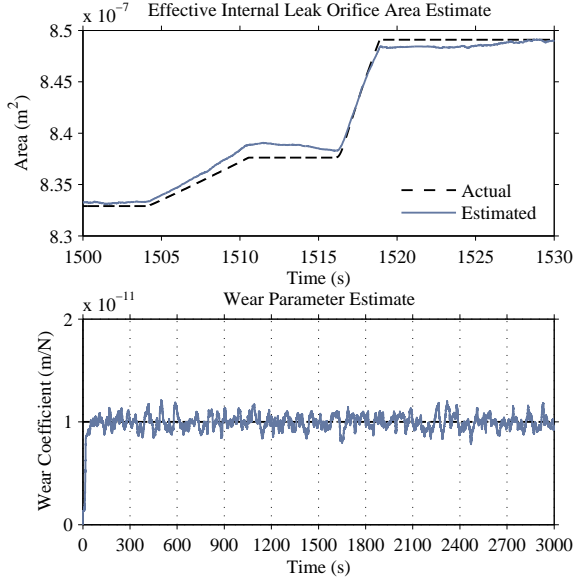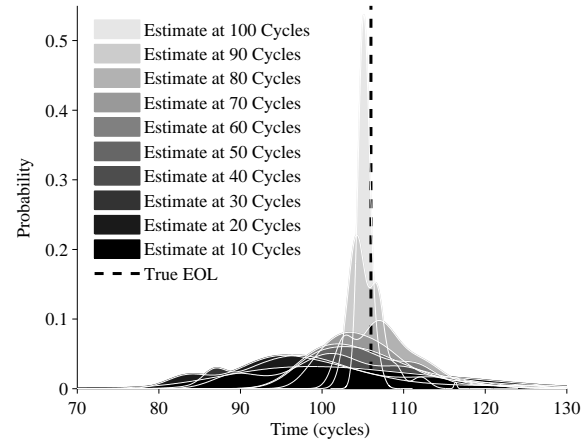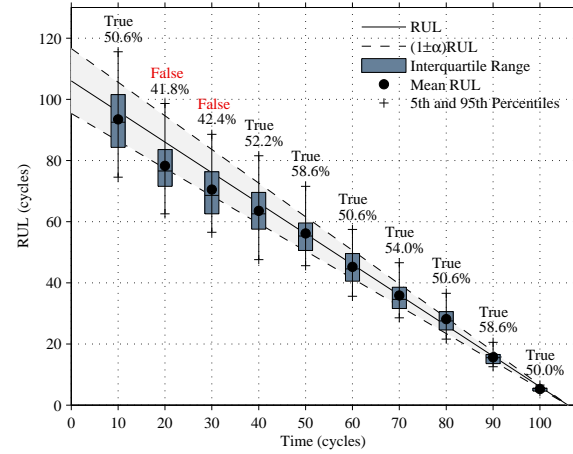
Figure 5: Estimation results for the growth of an internal leak, where $L = 3$.

Table 1: RUL predictions for internal leak

| $t_P$ | $RUL^*$ | $\overline{RUL}$ | $RA_{\text{Mean}}$ | $\widetilde{RUL}$ | $RA_{\text{Median}}$ | MAD |
|-------|---------|---------|---------|---------|---------|------|
| 10 | 96 | 93.45 | 0.973 | 92.52 | 0.964 | 8.79 |
| 20 | 86 | 78.23 | 0.910 | 76.58 | 0.890 | 6.02 |
| 30 | 76 | 70.51 | 0.928 | 68.58 | 0.902 | 6.04 |
| 40 | 66 | 63.54 | 0.963 | 62.57 | 0.948 | 6.05 |
| 50 | 56 | 56.13 | 0.998 | 55.28 | 0.987 | 4.70 |
| 60 | 46 | 45.20 | 0.983 | 44.52 | 0.968 | 4.86 |
| 70 | 36 | 35.85 | 0.996 | 34.59 | 0.961 | 4.26 |
| 80 | 26 | 28.12 | 0.919 | 27.53 | 0.941 | 3.00 |
| 90 | 16 | 15.63 | 0.977 | 15.52 | 0.970 | 1.28 |
| 100 | 6 | 5.26 | 0.876 | 5.47 | 0.911 | 0.88 |

version of the metric, which incorporates a third parameter, $\beta \in [0, 1]$, which defines a bound on the fraction of the probability mass of a prediction that falls within the $\alpha$-bounds (Saxena *et al.*, 2009). The metric evaluates to true at a given prediction point (i.e., a given $\lambda$) if the fraction of the probability mass within the $\alpha$-bounds, which we denote by $\pi_\alpha$, is greater than $\beta$. The metric, therefore, simultaneously captures aspects of both accuracy and precision. Fig. 7 shows the RUL predictions as box plots at each prediction point. The percent of the probability mass which falls within the $\alpha$-bounds is shown above each box plot, along with the outcome of the metric. The $\alpha$ and $\beta$ values would, in reality, be imposed by desired performance criteria of the prognostic system. Here, we choose reasonable values of $\alpha = 0.1$ and $\beta = 0.5$. In this case, the metric fails at the second and third prediction points, as less than half of the probability mass is contained within the $\alpha$-bounds. It should be noted, however, that the means of the distributions do fall within the bounds at those two points. If $\alpha$ is increased to $0.122$, the metric is satisfied at all points.



Figure 6: EOL predictions for the growth of an internal leak, where $L = 3$.



Figure 7: $\alpha$-$\lambda$ metric for the growth of an internal leak, where $L = 3$, $\alpha = 0.1$, and $\beta = 0.5$.

The predictions are quantified in Table 1, which provides the means and medians (in cycles) of the RUL distributions at each prediction point. Relative accuracy (RA) is also given, which, for a given prediction time $t_P$, is defined as (Saxena *et al.*, 2008; 2009):

$$RA_M(t_P) = 1 - \frac{|RUL^*(t_P) - M(RUL(t_P))|}{RUL^*(t_P)},$$

where $M$ denotes a selected measure of central tendency for the prediction distribution, and $RUL^*(t_P)$ denotes the true RUL at $t_P$. We use $\overline{RUL}$ to denote the mean of the distribution and $\widetilde{RUL}$ to denote the median, and $RA_{\text{Mean}}$ and $RA_{\text{Median}}$ to denote the relative accuracies computed with the mean and median, respectively, as the measures of central tendency. The table shows that RA is, on average, fairly high, and in this case, the mean provides a more accurate point estimate, as the average RA based on the mean is 0.952,

Table 2: Average performance for external leak predictions

| Fault | Lag | RMSE | $C_{w_{e,b}}$ | $\overline{RA}_{\text{Mean}}$ | $RA^-_{\text{Mean}}$ | $\overline{RA}_{\text{Median}}$ | $RA^-_{\text{Median}}$ | $\bar{\pi}_\alpha$ | $\pi^-_\alpha$ | MAD |
|---|---|---|---|---|---|---|---|---|---|---|
| Bottom External Leak | 0 | $7.74 \times 10^{-11}$ | 313.31 | 0.895 | 0.662 | 0.923 | 0.791 | 0.352 | 0.230 | 21.14 |
| ($w_{e,b} = 1 \times 10^{-9}$) | 1 | $7.42 \times 10^{-11}$ | 297.00 | 0.900 | 0.706 | 0.928 | 0.792 | 0.353 | 0.234 | 19.13 |
| | 2 | $7.00 \times 10^{-11}$ | 291.07 | 0.907 | 0.727 | 0.931 | 0.824 | 0.355 | 0.243 | 18.79 |
| | 3 | $7.19 \times 10^{-11}$ | 286.84 | 0.904 | 0.690 | 0.932 | 0.806 | 0.355 | 0.240 | 18.16 |
| | 4 | $7.67 \times 10^{-11}$ | 304.20 | 0.902 | 0.702 | 0.926 | 0.800 | 0.348 | 0.237 | 19.72 |
| Top External Leak | 0 | $4.61 \times 10^{-10}$ | 332.39 | 0.961 | 0.894 | 0.967 | 0.917 | 0.528 | 0.418 | 10.83 |
| ($w_{e,t} = 1 \times 10^{-8}$) | 1 | $4.68 \times 10^{-10}$ | 301.03 | 0.960 | 0.889 | 0.967 | 0.908 | 0.527 | 0.409 | 11.64 |
| | 2 | $4.88 \times 10^{-10}$ | 332.37 | 0.960 | 0.888 | 0.966 | 0.911 | 0.534 | 0.413 | 10.53 |
| | 3 | $4.91 \times 10^{-10}$ | 329.72 | 0.959 | 0.889 | 0.966 | 0.910 | 0.527 | 0.412 | 10.88 |
| | 4 | $4.86 \times 10^{-10}$ | 312.19 | 0.961 | 0.897 | 0.967 | 0.915 | 0.531 | 0.422 | 11.41 |

whereas based on the median, the average RA is $0.944$. The table also includes the median absolute deviation (MAD) of the distribution (in cycles) as a measure of variability. As EOL is approached, the MAD converges toward zero.

**Evaluation of Fixed-lag Filters**

We now examine how the use of fixed-lag filters can improve estimation and, subsequently, prediction. To illustrate this, we show results that quantify both estimation performance and prognostics performance, and demonstrate the link between the two aspects. Table 2 shows results for a comprehensive set of simulation experiments for a bottom external leak fault and a top external leak fault. For each lag, 15 experiments were performed. Each result in the table represents the particular metric averaged over those 15 experiments. Note that the same amount of noise was used in all experiments, which in this case was ten times the amount used in the internal leak example. Also, the wear coefficient for the top leak was chosen to be an order of magnitude larger than that for the bottom external leak, as this is the case to obtain similar EOL values (148 cycles for the bottom external leak, and 134 cycles for the top external leak).

We quantify estimation performance by the root mean square error (RMSE) of the weighted mean of the hidden parameter estimate from its true value. A smaller RMSE will entail a more accurate weighted mean of the RUL distribution, and, therefore, higher relative accuracy is achieved, and tighter $\alpha$-bounds can be met. In this case, for the bottom external leak, using a lag improves performance in estimation, and this corresponds also to improvements in prediction, as quantified by RA. A lag of 2 is optimal here, and as higher lags are considered the RMSE begins to increase. This occurs because the lookahead step becomes less reliable in noisy environments with higher lags. Process noise accumulates and predictions may lose their accuracy. However, the $L = 4$ case still outperforms the $L = 0$ case for RMSE. For the top external leak, the RMSE actually gets worse with a lag, but by a relatively small amount. In the context of prediction, this difference has no virtually no effect, as the remaining metrics have almost no difference between the cases. We attribute this to the fact that the wear parameter is an order of magnitude larger than for the bottom external leak, and at this scale, the effects of the fault are

easily distinguishable from noise even without lookahead.

We also report on the convergence of the damage parameter estimation error, computed as the distance from the origin to the centroid of the area under the error curve (Saxena *et al.*, 2008). A lower convergence score corresponds to faster convergence. Since the wear parameters are on the order of $10^{-9}$ and $10^{-8}$, the units of the convergence score are roughly in seconds. For the bottom external leak, the fixed lag cases clearly have better convergence, with $L = 3$ being optimal. Intuitively, this makes sense, as at the beginning of estimation, the lookahead allows particles farther from the initial guess, but closer to the true value, to be weighted more heavily, since $L$ steps ahead, these particles will be much more consistent with the observations than at the current time. A better convergence also means that a reliable prediction can be achieved earlier in time. For the top external leak, convergence scores are mostly the same, except with the $L = 1$ and $L = 4$ cases outperforming the other cases by approximately 30 and 20 seconds, respectively.

The average RA scores in the table represent the RA averaged over each prediction point of a single run (denoted as $\overline{RA}$), and the single number reported in the table for a given lag is this value averaged over the 15 experiments. We report RA calculated using both the mean and the median of the distributions, and here, the median provides a more accurate point estimate. We also report the average minimum RA, i.e., for a single run we take the minimum RA over all prediction points, and this number is averaged over the 15 experiments. This worst-case scenario corresponds to minimum $\alpha$-bounds that can always be satisfied. For the bottom external leak, the comparison between different lags is more pronounced. For the top leak, there is virtually no difference.

It is also important to examine the improvement in precision. This is quantified by the fraction of the prediction distributions that are within the $\alpha$-bounds (chosen as $0.1$). For a single run, we computed the average fraction over all prediction points, denoted as $\bar{\pi}_\alpha$. For a given lag, the table reports this value averaged over the 15 runs. We also report average minimum $\pi_\alpha$, denoted as $\pi^-_\alpha$, computed in the same way as for average minimum RA. Here, the difference is small. For the bottom external leak, when looking at the average

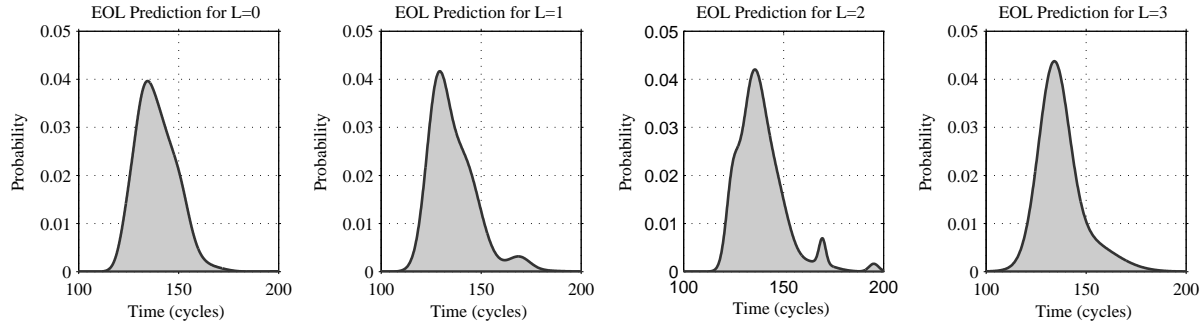Figure 8: EOL predictions for a top external leak with different lags.

Table 3: $\pi_\alpha^-$ values for different choices of $\alpha$ for a top external leak

| Lag | $\alpha = 0.15$ | $\alpha = 0.10$ | $\alpha = 0.05$ |
|---|---|---|---|
| 0 | $\pi_\alpha^- = 0.716$ | $\pi_\alpha^- = 0.478$ | $\pi_\alpha^- = 0.242$ |
| 1 | $\pi_\alpha^- = 0.770$ | $\pi_\alpha^- = 0.572$ | $\pi_\alpha^- = 0.294$ |
| 2 | $\pi_\alpha^- = 0.766$ | $\pi_\alpha^- = 0.556$ | $\pi_\alpha^- = 0.314$ |
| 3 | $\pi_\alpha^- = 0.788$ | $\pi_\alpha^- = 0.574$ | $\pi_\alpha^- = 0.290$ |
| 4 | $\pi_\alpha^- = 0.768$ | $\pi_\alpha^- = 0.578$ | $\pi_\alpha^- = 0.300$ |

worst case ($\pi_\alpha^-$), however, the difference is again more pronounced, meaning that for a given $\alpha$, the $\alpha$-$\lambda$ metric can be satisfied with larger $\beta$ values for fixed-lag filters. For the top leak, there is again little difference. We also report the MAD (in cycles), averaged over each run at the 30-cycle prediction point (the true RUL at this point is 118 cycles). An improvement in the fixed lag case is also visible here for the bottom external leak.

To further illustrate precision improvements, we show EOL predictions for the top external leak case with $L = 0$ to 3 in Fig. 8, under ten times less noise than the results reported in Table 2. EOL predictions were made at 50 cycles, where the true EOL is 134 cycles. The figure demonstrates that as $L$ is increased, the predictions become consistently more confident and, hence, more useful in making maintenance decisions. Table 3 shows this quantitatively, providing $\pi_\alpha^-$ values for different choices of $\alpha$. As the $\alpha$-bounds become tighter, the advantage of fixed-lag filters becomes more clear, as we see that a more significant portion of the probability mass is contained within tighter bounds as compared to the case with no lag. This also suggests that fixed-lag filters can have a more significant impact on performance when noise is small.

Overall, the results presented here show that fixed-lag filters can result in improved prognostics performance. While there is a visible performance increase in some cases, there is a trade-off between the extra computation involved in fixed-lag filters and the achieved performance gains. Compared to the case with no lag, a fixed-lag filter with lag $L$ will be doing $L + 1$ times the computations. If the performance requirement demands, for example, that average relative accuracy be 0.9, all values for $L$ presented in Table 2 satisfy this requirement, if $RA_{\text{Median}}$ is used, so, in this particular case, the extra computations offer no benefit

relative to the performance requirement. A fixed-lag approach may not be able to satisfy significantly different $\alpha$-bounds in which the $L = 0$ case cannot.

The choice of $L$ is also an important factor. In Table 2, $L = 2$ is optimal for the majority of the metrics for the bottom external leak fault. However, this may change depending on factors such as noise and sample time. The approximation to the smoothing distribution used here becomes less accurate as $L$ increases, so at some choice of $L$, it is no longer beneficial to be using a lag, i.e., estimation may become worse than the $L = 0$ case. The sample time also affects the choice of $L$, as with a smaller sample time, model uncertainty has less of an effect than with larger values of $L$, in which the uncertainty can build up over the lookahead window. The magnitude of the damage parameter also has an effect, as observed when comparing the results for the top and bottom external leaks. If the damage is progressing very slowly, then the progression is more distinguishable from noise if considered over a lag, so fixed-lag filters will be more beneficial. In practice, wear parameters would be much less than those considered here, so fixed-lag filters may offer more significant benefits than demonstrated here.[1]

## 6 CONCLUSIONS

In this paper, we presented a general model-based prognostics methodology using particle filters, formulated as a joint state-parameter estimation problem. State-parameter estimates are propagated forward in time to obtain EOL and RUL predictions, based on models that capture the progression of damage over time, characterized by a set of unknown parameters. We evaluated the use of fixed-lag particle filters within our scheme, and, overall, fixed-lag filters were shown to be able to provide improvements in state estimates and predictions, however, the improvements they offer depend on many factors, including the amount of noise and the magnitude of the wear parameters. Additional experiments are needed to investigate these issues more closely.

In this paper, we considered only single damage mechanisms active at any one time. This is not generally true, so, in future work, we will investigate the general case. It is well-known that standard particle

---

[1]Values of damage parameters were selected to so that EOL would be reached within reasonable experiment times.

filters become less effective at estimating multiple parameters, and additional techniques must be used, such as Rao-Blackwellization (e.g., see (Li *et al.*, 2007) for an application). In many cases, it is also safe to assume that noise is Gaussian, in which further improvements over the standard particle filter can be achieved. In future work we would like to investigate these cases and the performance gains that can be achieved, especially in the fixed-lag case.

**REFERENCES**

(Abbas *et al.*, 2007) M. Abbas, A. A. Ferri, M. E. Orchard, and G. J. Vachtsevanos. An intelligent diagnostic/prognostic framework for automotive electrical systems. In *2007 IEEE Intelligent Vehicles Symposium*, pages 352–357, 2007.

(Arulampalam *et al.*, 2002) M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp. A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE Transactions on Signal Processing*, 50(2):174–188, 2002.

(Byington *et al.*, 2004) C. S. Byington, M. Watson, D. Edwards, and P. Stoelting. A model-based approach to prognostics and health management for flight control actuators. In *Proceedings of the 2004 IEEE Aerospace Conference*, volume 6, pages 3551–3562, March 2004.

(Cappe *et al.*, 2007) O. Cappe, S. J. Godsill, and E. Moulines. An overview of existing methods and recent advances in sequential Monte Carlo. *Proceedings of the IEEE*, 95(5):899, 2007.

(Clapp and Godsill, 1999) T. C. Clapp and S. J. Godsill. Fixed-lag smoothing using sequential importance sampling. *Bayesian Statistics IV*, pages 743–752, 1999.

(Doucet *et al.*, 2000) A. Doucet, S. Godsill, and C. Andrieu. On sequential Monte Carlo sampling methods for Bayesian filtering. *Statistics and Computing*, 10:197–208, 2000.

(Hutchings, 1992) I. M. Hutchings. *Tribology: friction and wear of engineering materials*. Hodder & Stoughton Publishers, 1992.

(Kitagawa, 1996) G. Kitagawa. Monte Carlo filter and smoother for non-Gaussian nonlinear state space models. *Journal of Computational and Graphical Statistics*, 5(1):1–25, 1996.

(Li *et al.*, 2007) P. Li, R. Goodall, P. Weston, C. Seng Ling, C. Goodman, and C. Roberts. Estimation of railway vehicle suspension parameters for condition monitoring. *Control Engineering Practice*, 15(1):43–55, 2007.

(Liu and West, 2001) J. Liu and M. West. Combined parameter and state estimation in simulation-based filtering. *Sequential Monte Carlo methods in Practice*, pages 197–223, 2001.

(Orchard *et al.*, 2008) M. Orchard, G. Kacprzynski, K. Goebel, B. Saha, and G. Vachtsevanos. Advances in uncertainty representation and management for particle filtering applied to prognostics. In *Proceedings of International Conference on Prognostics and Health Management*, Oct 2008.

(Perry and Green, 2007) R.H. Perry and D.W. Green. *Perry's chemical engineers' handbook*. McGraw-Hill Professional, 2007.

(Roemer *et al.*, 2005) M. Roemer, C. Byington, G. Kacprzynski, and G. Vachtsevanos. An overview of selected prognostic technologies with reference to an integrated PHM architecture. In *Proceedings of the First International Forum on Integrated System Health Engineering and Management in Aerospace*, 2005.

(Saha and Goebel, 2009) B. Saha and K. Goebel. Modeling Li-ion battery capacity depletion in a particle filtering framework. In *Proceedings of the Annual Conference of the Prognostics and Health Management Society 2009*, September 2009.

(Saxena *et al.*, 2008) A. Saxena, J. Celaya, E. Balaban, K. Goebel, B. Saha, S. Saha, and M. Schwabacher. Metrics for evaluating performance of prognostic techniques. In *International Conference on Prognostics and Health Management 2008*, Oct 2008.

(Saxena *et al.*, 2009) A. Saxena, J. Celaya, B. Saha, S. Saha, and K. Goebel. On applying the prognostic performance metrics. In *Proceedings of the Annual Conference of the Prognostics and Health Management Society 2009*, September 2009.

(Schwabacher, 2005) M.A. Schwabacher. A survey of data-driven prognostics. In *Proceedings of the AIAA Infotech@ Aerospace Conference*, 2005.