

An Integrated Model-Based Diagnostic and Prognostic Framework

Indranil Roychoudhury¹ and Matthew Daigle²

¹ SGT Inc., NASA Ames Research Center, Moffett Field, CA, 94035, USA
indranil.roychoudhury@nasa.gov

² University of California, Santa Cruz, NASA Ames Research Center, Moffett Field, CA, 94035, USA
matthew.j.daigle@nasa.gov

ABSTRACT

Systems health monitoring is essential in guaranteeing the safe, efficient, and correct operation of complex engineered systems. Diagnosis, which consists of detection, isolation and identification of faults, and prognosis, which consists of prediction of the remaining useful life of components, subsystems, or systems constitute system health monitoring. This paper presents a comprehensive framework for integrated model-based diagnosis and prognosis of complex systems, where we make use of a common modeling paradigm to model both the nominal and faulty behavior in all aspects of health monitoring. We illustrate our approach using a simulated propellant loading system that includes tanks, valves, and pumps.

1 INTRODUCTION

Systems health monitoring is essential in guaranteeing the safe, efficient, and correct operation of complex engineered systems. The integral tasks of systems health monitoring include diagnostics and prognostics. Diagnosis involves *detecting* when a fault has occurred, *isolating* the true fault from many possible fault candidates, and *identifying* the true damage to the system. While diagnosis involves determining what *has happened* to the system, prognosis, on the other hand, involves determining what *will happen*. Specifically, prognosis involves *predicting* how much useful life remains in the different components, subsystems, or systems. Based on these predictions, effective actions can be taken to minimize (or completely remove) any loss of life or property, optimize maintenance, and extend component life.

A large body of research exists for both diagnostics and prognostics. However, many diagnosis approaches stop at the fault isolation step, and seldom perform fault identification, and most prognostic approaches assume some diagnosis has been performed and focus on prognosis of a single failure mode. This paper presents an integrated framework for model-based diagnostics and prognostics of complex systems, in

which we make use of a common modeling framework for modeling both the nominal and faulty system behavior. We assume only single faults in this paper.

In our approach, we start with modeling the nominal system behavior, as well as how different faults manifest in the system behavior and progress over time. An observer built with the nominal model is used to generate estimates of nominal system behavior, and when the deviation of observed measurements from the nominal estimates is statistically significant, a fault is detected. Fault isolation involves comparing the observed measurement deviations to predictions of how these measurements would deviate for different possible faults, and removing from consideration fault candidates that are inconsistent with the observed deviations. Fault identification involves tracking the observed system measurements using multiple observers, each built with a hypothesized fault model integrated with the nominal model, and performing joint state-parameter estimation (Roychoudhury, 2009). The prognosis module predicts the remaining useful life of a component, subsystem, or system, using, for each hypothesized fault, a predictor based on the fault progression model integrated with the nominal model (Daigle and Goebel, 2010). We demonstrate our approach on a simulated propellant loading system. Our experiments illustrate that our integrated diagnostic and prognostic approach diagnosed faults and predicted the EOL accurately.

This paper is organized as follows. Section 2 provides the problem formulation and architecture for our diagnostic and prognostic framework, and Section 3 describes the different components of our integrated diagnostic and prognostic approach. The case study and experimental results are presented in Section 4. Section 5 concludes the paper.

2 DIAGNOSTIC AND PROGNOSTIC ARCHITECTURE

In this section, we formulate the diagnosis and prognosis problem, and provide an architecture for an integrated diagnosis and prognosis approach.

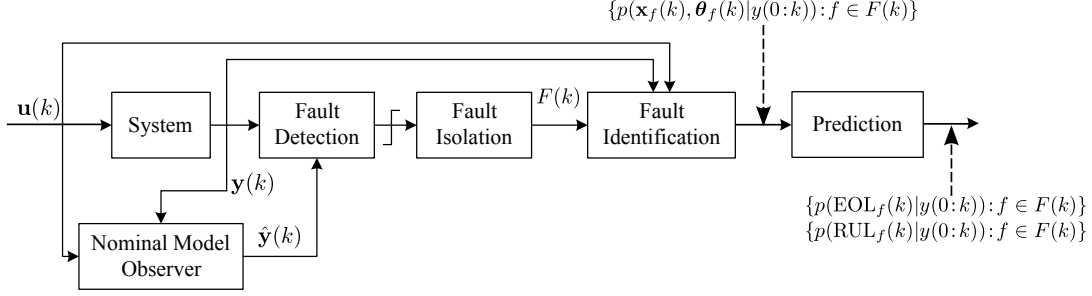


Figure 1: The integrated diagnostic and prognostic architecture.

2.1 Problem Formulation

We define a system model for representing system behavior under nominal operation, as follows:

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{f}(t, \mathbf{x}(t), \boldsymbol{\theta}(t), \mathbf{u}(t), \mathbf{v}(t)) \\ \mathbf{y}(t) &= \mathbf{h}(t, \mathbf{x}(t), \boldsymbol{\theta}(t), \mathbf{u}(t), \mathbf{n}(t)),\end{aligned}$$

where $\mathbf{x}(t) \in \mathbb{R}^{n_x}$ is the state vector, $\boldsymbol{\theta}(t) \in \mathbb{R}^{n_\theta}$ is the parameter vector, $\mathbf{u}(t) \in \mathbb{R}^{n_u}$ is the input vector, $\mathbf{v}(t) \in \mathbb{R}^{n_v}$ is the process noise vector, \mathbf{f} is the state equation, $\mathbf{y}(t) \in \mathbb{R}^{n_y}$ is the output vector, $\mathbf{n}(t) \in \mathbb{R}^{n_n}$ is the measurement noise vector, and \mathbf{h} is the output equation. The parameters $\boldsymbol{\theta}(t)$ evolve in an unknown way, but, in practice, are usually considered as constants.

Any change in the above nominal system model represents a fault. In this work, we restrict faults solely to changes in system parameters, $\boldsymbol{\theta}(t)$. Under the single fault assumption, only one parameter can deviate from nominal. Hence, we denote a fault, $f \in F$, as a tuple, (θ, g_f) , where, $\theta \in \boldsymbol{\theta}$ is the faulty parameter, and g_f denotes the *fault progression function* according to which, fault f is manifested in parameter θ , i.e.,

$$\dot{\theta}(t) = g_f(t, \mathbf{x}_f(t), \boldsymbol{\theta}_f(t), \mathbf{u}(t), \mathbf{m}_f(t)),$$

where $\mathbf{x}_f(t) = [\mathbf{x}(t), \theta(t)]^T$, $\boldsymbol{\theta}_f(t) = [\boldsymbol{\theta}(t) \setminus \{\theta(t)\}, \phi_f(t)]^T$, $\phi_f(t) \in \mathbb{R}^{n_{\phi_f}}$ is a vector of fault progression model parameters, and $\mathbf{m}_f(t) \in \mathbb{R}^{n_{m_f}}$ is a process noise vector associated with the fault progression model.

The single fault assumption also implies that the faulty system model for fault $f = (\theta, g_f)$ involves integrating a single fault progression model into the nominal system model described above, as shown below:

$$\begin{aligned}\dot{\mathbf{x}}_f(t) &= \mathbf{f}_f(t, \mathbf{x}_f(t), \boldsymbol{\theta}_f(t), \mathbf{u}(t), \mathbf{v}(t)) \\ \mathbf{y}(t) &= \mathbf{h}(t, \mathbf{x}(t), \boldsymbol{\theta}(t), \mathbf{u}(t), \mathbf{n}(t)),\end{aligned}$$

where,

$$\mathbf{f}_f(\cdot) = \begin{bmatrix} \mathbf{f}(t, \mathbf{x}(t), \boldsymbol{\theta}(t), \mathbf{u}(t), \mathbf{v}(t)) \\ g_f(t, \mathbf{x}_f(t), \boldsymbol{\theta}_f(t), \mathbf{u}(t), \mathbf{m}(t)) \end{bmatrix} = \begin{bmatrix} \dot{\mathbf{x}}(t) \\ \dot{\theta}(t) \end{bmatrix}.$$

Since any of the several parameters in a system model can be faulty, the goal of diagnosis is to:

1. Detect a change in some $\theta \in \boldsymbol{\theta}$;
2. Isolate, under the single fault assumption, the true $f \in F$, i.e., both the parameter θ that has changed, and its fault progression model g_f ; and

3. Identify the extent of damage by computing $p(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k) | \mathbf{y}(0:k))$, where $\mathbf{y}(0:k)$ denotes all measurements observed up to the present discrete time step, k , and $\mathbf{x}_f(k)$ and $\boldsymbol{\theta}_f(k)$ denote the value of \mathbf{x}_f and $\boldsymbol{\theta}_f$ at time step k , respectively.

The goal of prognosis is to predict for a given fault, f , using, $p(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k) | \mathbf{y}(0:k_P))$, a probability distribution of end of life (EOL), i.e., $p(\text{EOL}_f(k_P) | \mathbf{y}(0:k_P))$, and/or remaining useful life (RUL), i.e., $p(\text{RUL}_f(k_P) | \mathbf{y}(0:k_P))$ at a given time point k_P (Daigle and Goebel, 2010). We predict the probability distribution, rather than the exact EOL and/or RUL, since, there is inherent uncertainty in the prediction of state, as well as, the future input uncertainty. A set of constraints define the acceptable behavior of a system. EOL is reached when one or more of the constraints are no more met. We define $T_{\text{EOL}_f} = 1$ if these constraints are valid, and $T_{\text{EOL}_f} = 0$ otherwise.

So, EOL_f may be defined as $\text{EOL}_f(k_P) \triangleq \inf\{k \in \mathbb{R} : k \geq k_P \text{ and } T_{\text{EOL}_f}(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k)) = 1\}$. i.e., EOL is the earliest time point at which the threshold is reached. Given $\text{EOL}_f(k_P)$, RUL may then be defined with $\text{RUL}_f(k_P) \triangleq \text{EOL}_f(k_P) - k_P$.

2.2 Architecture

Fig. 1 illustrates the architecture of our combined diagnostic and prognostic scheme. At each discrete time, k , the system takes as inputs $\mathbf{u}(k)$, and outputs measurements $\mathbf{y}(k)$. The nominal observer also takes as inputs $\mathbf{u}(k)$, and generates estimates of nominal measurements, $\hat{\mathbf{y}}(k)$. The fault detector then takes in the observed and estimated measurements, $\mathbf{y}(k)$ and $\hat{\mathbf{y}}(k)$, and detects when a fault has occurred based on the residual, $\mathbf{r}(k) = \mathbf{y}(k) - \hat{\mathbf{y}}(k)$. Once a fault is detected, fault isolation is initiated. The fault isolation block takes as inputs the measurement residuals. These measurement residuals are used along with predictions of how each measurement is expected to deviate from nominal for each possible fault in the system to generate a set of fault candidates $F(k)$ at time k that explain the observed deviations in measurements till time k . The fault identification module, for each fault, $f \in F(k)$, estimates $p(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k) | \mathbf{y}(0:k))$. Finally, the prediction module takes in as input $p(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k) | \mathbf{y}(0:k))$ to make predictions of EOL and/or RUL, i.e., $p(\text{EOL}_f(k) | \mathbf{y}(0:k))$ and/or $p(\text{RUL}_f(k) | \mathbf{y}(0:k))$.

3 DIAGNOSIS AND PROGNOSIS APPROACH

In this section, we describe how we implement each of the different modules of the integrated diagnosis and prognosis architecture.

3.1 Nominal Observer

The nominal observer takes as inputs the system inputs, $\mathbf{u}(k)$, and measurements, $\mathbf{y}(0 : k)$, and the initial state of the system, and uses the state transition function, $\mathbf{f}(\cdot)$ and observation function $\mathbf{h}(\cdot)$ to estimate distributions of states, $\mathbf{x}(k)$ and parameters, $\boldsymbol{\theta}(k)$, i.e., $p(\mathbf{x}(k), \boldsymbol{\theta}(k) | \mathbf{y}(0:k))$.

While any standard filtering scheme, e.g., Kalman filter, extended Kalman filter, unscented Kalman filter, among others, can be adopted as our nominal observer, we adopt the particle filter as a general solution (Arunlampalam *et al.*, 2002). Particle filtering is the most general estimation scheme as it can be applied to nonlinear systems with arbitrary probability distributions for measurement noise and modeling error that can be nonlinearly coupled with the states. Particle filtering is a sequential Monte Carlo sampling method for Bayesian filtering and approximates the belief state of a system using a weighted set of samples, or particles. Each sample, or particle, consists of an instantiation of values of the state vector, and describes a possible system state. As more observations are obtained, each particle is moved stochastically to a new state using the nominal state transition function, and the weight of each particle is readjusted to reflect the likelihood of that observation given the particle's new state.

3.2 Fault Detection

A fault is detected when a residual, $r(k) \in \mathbf{r}(k)$, i.e., the difference between the observed (faulty) and estimated (nominal) values of a measurement is determined to be statistically significant (Daigle *et al.*, 2010). In our work, we use a Z -test coupled with a sliding window technique to determine this statistical significance. Our fault detection scheme is described in detail in (Daigle *et al.*, 2010).

3.3 Fault Isolation

Once a fault is detected, at each subsequent time step, every measurement residual is qualitatively abstracted into a tuple of qualitative symbols, (σ_1, σ_2) , where $\sigma_1 \in \{0, +, -\}$ represents the qualitative magnitude change, and $\sigma_2 \in \{0, +, -\}$ represents the qualitative slope change. The symbols, 0, +, or -, denote whether the magnitude or slope of this measurement is at, above, or below nominal, respectively. The symbols are generated using a sliding window technique as described in detail in (Mosterman and Biswas, 1999).

Based on the first observed statistically significant measurement deviation, we first generate a set of possible fault candidates. Then, for each fault candidate, we determine a fault signature for each measurement. A fault signature of a fault for a measurement is a prediction of how the measurement will deviate from nominal under the effect of the particular fault. Fault signatures are also of the form (s_1, s_2) , where $s_1 \in \{0, +, -\}$ and $s_2 \in \{0, +, -\}$ capture qualitatively the direction of change to be expected in the

magnitude and slope of each measurement from nominal if the fault occurs. The procedure for generating fault signatures is presented in detail in (Mosterman and Biswas, 1999).

Given the set of fault candidates and fault signatures, as more measurements deviate from nominal, the fault signatures are compared to the observed measurement deviations (captured symbolically) and any fault candidate whose fault signature is inconsistent with the observed measurement deviation is removed from consideration. As more and more measurements are observed to deviate from nominal, the fault candidate set will reduce, ideally resulting in a singleton (since we assume single faults).

However, in some cases, the qualitative fault signatures alone are not sufficient in distinguishing all faults, or fault effects may take too long to manifest, and quantitative analysis is needed to correctly diagnose the true fault.

3.4 Fault Identification

We initiate the quantitative fault identification procedure after the qualitative fault signature-based isolation scheme is executed for p time steps or till the number of fault candidates reduces to less than s , whichever is achieved first. The design parameters p and s are chosen empirically for each domain and system on a case-by-case basis. If multiple fault candidates are valid when fault identification is initiated, fault identification also helps in reduction of inconsistent fault candidates, as described next.

Once fault identification is invoked, under the single fault assumption, for each remaining fault candidate, f , we instantiate a particle filter-based observer using its faulty system model, \mathbf{f}_f , generated, as described in Section 2.1, by extending the nominal system model with the fault progression model. Then each fault observer tracks the observed system measurements independently, and generates estimates of $\hat{\mathbf{y}}(k)$ and $p(\mathbf{x}_f(k), \boldsymbol{\theta}_f(k) | \mathbf{y}(k_d - \Delta k^{max} : k))$, where $\Delta k^{max} \geq k_d - k_f$ is the maximum delay assumed to be possible between the time of fault occurrence, k_f , and the time of fault detection, k_d . Each fault observer is initialized to known or estimated values of \mathbf{x} and $\boldsymbol{\theta}$ in the nominal operating region, and ϕ_f is initialized to a zero vector. If multiple fault candidates remain when fault identification is invoked, for each fault observer, a Z -test is used to determine if the deviation of a measurement estimated by the particle filter from the corresponding actual observation is statistically significant. Since we are considering only single faults, the expectation is that eventually, the estimates of only the correct fault observer will converge to the observed measurements, while those of all others will deviate from the observed measurements. Thus fault identification also helps in fault isolation. Practically, even the true fault model will take some time before tracking the measurements correctly, since initially, the system and damage parameter values are guesses. We assume that the true fault observer will converge to the observed measurements within s_d time steps of its invocation. Thus, the Z -tests are monitored only after time s_d time steps are over. For details of our fault diagnosis approach, please refer to (Roychoudhury, 2009).

Algorithm 1 EOL Prediction

Inputs: $\{\mathbf{x}_f^i(k_P), \boldsymbol{\theta}_f^i(k_P), w^i(k_P)\}_{i=1}^N$
Outputs: $\{EOL_f^i(k_P), w^i(k_P)\}_{i=1}^N$
for $i = 1$ **to** N **do**
 $k \leftarrow t_P$
 $\mathbf{x}_f^i(k) \leftarrow \mathbf{x}_f^i(k_P)$
 $\boldsymbol{\theta}_f^i(k) \leftarrow \boldsymbol{\theta}_f^i(k_P)$
 while $T_{EOL_f}(\mathbf{x}_f^i(k), \boldsymbol{\theta}_f^i(k)) = 0$ **do**
 Predict $\hat{\mathbf{u}}(k)$
 $\boldsymbol{\theta}_f^i(k+1) \sim p(\boldsymbol{\theta}_f(k+1)|\boldsymbol{\theta}_f^i(k))$
 $\mathbf{x}_f^i(k+1) \sim p(\mathbf{x}_f(k+1)|\mathbf{x}_f^i(k), \boldsymbol{\theta}_f^i(k), \hat{\mathbf{u}}(k))$
 $k \leftarrow k+1$
 $\mathbf{x}_f^i(k) \leftarrow \mathbf{x}_f^i(k+1)$
 $\boldsymbol{\theta}_f^i(k) \leftarrow \boldsymbol{\theta}_f^i(k+1)$
 $EOL_f^i(k_P) \leftarrow k$

3.5 Prediction

The prediction module is invoked at time k_P to predict the EOL and/or RUL of the component for each hypothesized fault, f . Specifically, using the current joint state-parameter estimate, $p(\mathbf{x}_f(k_P), \boldsymbol{\theta}_f(k_P)|\mathbf{y}(0:k_P))$, which represents the most up-to-date knowledge of the system at time k_P , the goal is to compute $p(EOL_f(k_P)|\mathbf{y}(0:k_P))$ and $p(RUL_f(k_P)|\mathbf{y}(0:k_P))$. We assume the state-parameter distribution is represented as a discrete set of weighted samples, i.e.,

$$p(\mathbf{x}_f(k_P), \boldsymbol{\theta}_f(k_P)|\mathbf{y}(0:k_P)) \approx \sum_{i=1}^N w^i(k_P) \delta_{(\mathbf{x}_f^i(k_P), \boldsymbol{\theta}_f^i(k_P))} (d\mathbf{x}_f(k_P) d\boldsymbol{\theta}_f(k_P)),$$

where i denotes the index of a single sample, w^i is the weight of this sample, and δ represents the Dirac delta function.

Similarly, we can approximate the EOL as

$$p(EOL_f(k_P)|\mathbf{y}(0:k_P)) \approx \sum_{i=1}^N w^i(k_P) \delta_{EOL_f^i(k_P)} (dEOL_f(k_P)).$$

To compute EOL, then, we propagate each sample in state-parameter distribution forward to its own EOL and use that sample's weight at k_P for the weight of its EOL prediction.

The general approach to solving the prediction problem is through simulation. Each sample is simulated forward to EOL to obtain the complete EOL distribution. The pseudocode for the prediction procedure is given as Algorithm 1 (Daigle and Goebel, 2010). Each sample i is propagated forward until $T_{EOL_f}(\mathbf{x}_f^i(k), \boldsymbol{\theta}_f^i(k))$ evaluates to 1; at this point EOL has been reached for this particle. In this work, we adopt particle filter-based approach for prediction.

Note that we need to hypothesize future inputs of the system, $\hat{\mathbf{u}}(k)$, for prediction, since fault progression is dependent on the operational conditions of the system. The choice of expected future inputs depends on the knowledge of operational settings.

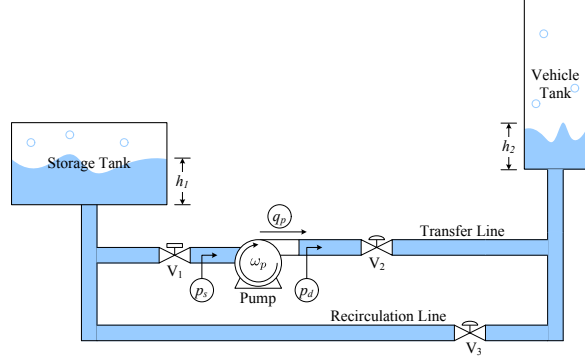


Figure 2: Fueling system schematic.

4 CASE STUDY

We apply the approach to a simulation of a fueling system. The system schematic is shown in Fig. 2 and is based on a subset of the system presented in (Goodrich *et al.*, 2009). Liquid is drained from a storage tank through a transfer line via a pump, into a vehicle tank. In normal operation, both valves V_1 and V_2 on the transfer line are fully open, and the valve V_3 on the recirculation line is fully closed.

Measurements include the tank heights, h_1 and h_2 , the suction and discharge pressures of the pump, p_s and p_d , the rotational velocity of the pump, ω_p , the discharge flow of the pump, q_p , and the thrust bearing, radial bearing, and oil temperatures of the pump, T_t , T_r , and T_o , respectively (the location of temperature sensors are not shown in Fig. 2).

In the remainder of this section, we first describe the system model. We then provide an example scenario to demonstrate the approach, followed by a summary of diagnosis and prognosis results.

4.1 System Modeling

The storage tank and vehicle tank masses are described by

$$\begin{aligned} \dot{m}_1(t) &= q_{V3} - q_{V1} - q_{l1} \\ \dot{m}_2(t) &= q_{V2} - q_{V3} - q_{l2}, \end{aligned}$$

respectively, where the flows are defined as

$$\begin{aligned} q_{V1} &= u_{V1} A_{V1} c_{V1} \sqrt{|p_1 - p_s|} \text{sign}(p_1 - p_s) \\ q_{V2} &= u_{V2} A_{V2} c_{V2} \sqrt{|p_d - p_2|} \text{sign}(p_d - p_2) \\ q_{V3} &= u_{V3} A_{V3} c_{V3} \sqrt{|p_2 - p_1|} \text{sign}(p_2 - p_1) \\ q_{l1} &= A_{l1} \sqrt{|p_1 - p_{atm}|} \text{sign}(p_1 - p_{atm}) \\ q_{l2} &= A_{l2} \sqrt{|p_2 - p_{atm}|} \text{sign}(p_2 - p_{atm}), \end{aligned}$$

such that $u_{Vi} \in [0, 1]$ denotes the commanded position of valve V_i with 0 denoting the valve is fully closed, and 1 denoting the valve is fully open; c_C denotes the capacitance of component C ; and A_C denotes the product of the cross-sectional area of component C and the flow coefficient of component C , c_C . The tank pressures are given by

$$\begin{aligned} p_1 &= p_{atm} + \rho g h_1 \\ p_2 &= p_{atm} + \rho g h_2, \end{aligned}$$

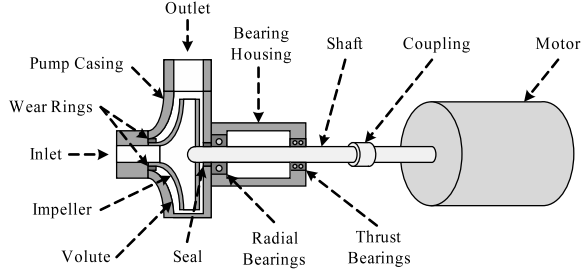


Figure 3: Centrifugal pump.

with $h_j = m_j/(\rho A_j)$, where ρ is the liquid density and A_j is the tank cross-sectional area (the tanks are assumed to have a uniform cross-sectional area). The suction and discharge pressures are given by

$$\begin{aligned}\dot{p}_s &= 1/C_s(q_{V1} - q_p) \\ \dot{p}_d &= 1/C_d(q_p - q_{p2}),\end{aligned}$$

where C_s and C_d are pipe capacitances, and q_p is the pump flow.

The centrifugal pump takes in fluid through its inlet, and the rotation of its impeller forces the fluid through the outlet. Fig. 3 presents the schematic of a centrifugal pump.

The rotational velocity of the pump is described using a torque balance,

$$\dot{\omega}_p = \frac{1}{J} (\tau_e(t) - r\omega_p(t) - \tau_L(t)),$$

where J is the lumped motor/pump inertia, τ_e is the electromagnetic torque provided by the motor, r is the lumped friction parameter, and τ_L is the load torque. A torque is produced on the rotor only when there is a difference (i.e., slip) between the synchronous speed of the supply voltage, ω_s and the mechanical rotation, ω_p , where slip s is defined as

$$s = \frac{\omega_s - \omega_p}{\omega_s}.$$

The expression for the torque τ_e for an alternating-current induction motor is (Lyshevski, 1999):

$$\tau_e = \frac{npR_2}{s\omega_s} \frac{V_{rms}^2}{(R_1 + R_2/s)^2 + (\omega_s L_1 + \omega_s L_2)^2},$$

where R_1 is the stator resistance, L_1 is the stator inductance, R_2 is the rotor resistance, L_2 is the rotor inductance, n is the number of phases (typically 3), and p is the number of magnetic pole pairs. The rotor speed may be controlled by changing the input frequency ω_s .

The load torque τ_L is a polynomial function of the flow rate through the pump and the impeller rotational velocity (Wolfram *et al.*, 2001; Kallesøe, 2005):

$$\tau_L = a_0\omega_p^2 + a_1\omega_p q_p - a_2 q_p^2,$$

where q_p is the pump flow, and a_0 , a_1 , and a_2 are coefficients derived from the pump geometry.

The rotation of the impeller creates a pressure difference from the inlet to the outlet of the pump, which

drives the pump flow, q_p . The pump pressure is computed as

$$p_p = b_0\omega_p^2 + b_1\omega_p q_p - b_2 q_p^2,$$

where b_0 , b_1 , and b_2 are coefficients derived from the pump geometry. The parameter b_0 is proportional to impeller area A_i . Flow through the impeller, q_i , is computed using the pressure differences:

$$q_i = c\sqrt{|p_s + p_p - p_d| \text{sign}(p_s + p_p - p_d)},$$

where c is a flow coefficient, p_s is the suction pressure, and p_d is the discharge pressure. The small (normal) leakage flow from the discharge end to the suction end due to the clearance between the wear rings and the impeller is described by

$$q_l = c_l\sqrt{|p_d - p_s| \text{sign}(p_d - p_s)},$$

where c_l is a flow coefficient. The discharge flow, q_p , is then

$$q_p = q_i - q_l.$$

Pump temperatures are often monitored as indicators of pump condition. The oil heats up due to the radial and thrust bearings and cools to the environment:

$$\begin{aligned}\dot{T}_o &= \frac{1}{J_o} (H_{o,1}(T_t - T_o) + H_{o,2}(T_r - T_o) - \\ &H_{o,3}(T_o - T_a)),\end{aligned}$$

where J_o is the thermal inertia of the oil, and the $H_{o,i}$ terms are heat transfer coefficients. The thrust bearings heat up due to the friction between the pump shaft and the bearings, and cool to the oil and the environment:

$$\dot{T}_t = \frac{1}{J_t} (r_t\omega^2 - H_{t,1}(T_t - T_o) - H_{t,2}(T_t - T_a)),$$

where J_t is the thermal inertia of the thrust bearings, r_t is the friction coefficient for the thrust bearings, and the $H_{t,i}$ terms are heat transfer coefficients. The radial bearings behave similarly:

$$\dot{T}_r = \frac{1}{J_r} (r_r\omega^2 - H_{r,1}(T_r - T_o) - H_{r,2}(T_r - T_a))$$

where J_r is the thermal inertia of the radial bearings, r_r is the friction coefficient for the radial bearings, and the $H_{r,i}$ terms are heat transfer coefficients. Note that r_t and r_r contribute to the overall friction coefficient r .

Faulty System Modeling

For our experiments, we consider the eight faults shown in Table 1. Either tank can have a leak fault, represented as abrupt increase in parameter A_{l1} or A_{l2} . The abrupt increase in A_{l1} is characterized by the fault progression function,

$$\dot{A}_{l1} = \begin{cases} \delta(t_f)\Delta A_{l1}, & t = t_f \\ 0, & \text{otherwise} \end{cases}$$

where δ represents a Dirac delta function, t_f denotes the time of fault occurrence, and ΔA_{l1} is the fault parameter. The fault progression function for leak fault A_{l2}^+ also takes a similar form.

Valves V_1 and V_2 are nominally open and valve V_3 is nominally closed. Hence, stuck faults in these three valves are denoted by x_1^- , x_2^- , and x_3^+ where Δx_1 , Δx_2 , or Δx_3 denote the values at which these valves get abruptly stuck at, respectively. Therefore, the fault progression function for valve V_1 is

$$\dot{x}_1 = \begin{cases} \delta(t_f)\Delta x_1, & t = t_f \\ 0, & \text{otherwise} \end{cases}$$

Faults x_2^- and x_3^+ have similar fault progression functions. For abrupt faults, the component is assumed to have reached its EOL, i.e., $T_{EOL_f} = 1$, as soon as the fault occurs, i.e., as soon as a leak is present in a tank, or a valve becomes stuck. As a result, RUL predictions associated with these components are trivially 0 whenever they are diagnosed.

However, faults in the pump are not abrupt, but incipient, i.e., they progress slowly. These pump faults include impeller wear, represented as a progressive change in impeller area A_i ; bearing wear, represented as progressive changes in the thrust bearing friction coefficient r_t or the radial bearing friction coefficient r_r .

The impeller wear, A_i^- , is represented by a gradual decrease in impeller area A_i (Biswas and Mahadevan, 2007; Tu *et al.*, 2007). Since the impeller area is proportional to b_0 , a decrease causes a decrease in the pump pressure, and, hence, the pump efficiency. The erosive wear equation (Hutchings, 1992) is used to describe the how the impeller area changes over time. The erosive wear rate is proportional to fluid velocity times friction force. Fluid velocity is proportional to volumetric flow rate, and friction force is proportional to fluid velocity. We lump the proportionality constants into the wear coefficient w_A to obtain the fault progression function for A_i^- as follows:

$$\dot{A}_i(t) = -w_{A_i} q_i(t)^2.$$

We represent the bearing wear faults, r_r^+ or r_t^+ , as gradual increases in the friction coefficients of sliding and rolling friction, r_t and r_r , respectively, due to wear of materials (Hutchings, 1992; Daigle and Goebel, 2010), and modeled as:

$$\begin{aligned} \dot{r}_t(t) &= w_t r_t \omega^2 \\ \dot{r}_r(t) &= w_r r_r \omega^2, \end{aligned}$$

where w_t and w_r are the wear coefficients. The slip compensation provided by the electromagnetic torque generation masks small changes in friction, so it is only with very large increases that a change in ω will be observed, but small changes produce easily observable changes in temperature.

The pump is still functional, i.e., it is still delivering fluid, in the presence the three wear faults. Hence, its EOL is defined by the effective impeller area decreasing to a certain level A_i^\downarrow , and by its temperatures exceeding given thresholds at which irreversible damage occurs, T_t^\uparrow , T_r^\uparrow , or T_o^\uparrow , where abnormal temperature increases are related to increases in bearing friction. So for a pump fault $f \in F$, $T_{EOL_f} = 1$ if $A_i(t) \leq A_i^\downarrow$, $T_t(t) \geq T_t^\uparrow$, $T_r(t) \geq T_r^\uparrow$, or $T_o(t) \geq T_o^\uparrow$.

Table 1: Faults of Interest

Fault Name	Description	θ	g_f	ϕ_f
A_{I1}^+	Leak in storage tank	A_{I1}	$\dot{A}_{I1} = \begin{cases} \delta(t_f)\Delta A_{I1}, & t = t_f \\ 0, & \text{otherwise} \end{cases}$	ΔA_{I1}
A_{I2}^+	Leak in vehicle tank	A_{I2}	$\dot{A}_{I2} = \begin{cases} \delta(t_f)\Delta A_{I2}, & t = t_f \\ 0, & \text{otherwise} \end{cases}$	ΔA_{I2}
x_1^-	V_1 stuck	x_1	$\dot{x}_1 = \begin{cases} \delta(t_f)\Delta x_1, & t = t_f \\ 0, & \text{otherwise} \end{cases}$	Δx_1
x_2^-	V_2 stuck	x_2	$\dot{x}_2 = \begin{cases} \delta(t_f)\Delta x_2, & t = t_f \\ 0, & \text{otherwise} \end{cases}$	Δx_2
x_3^+	V_3 stuck	x_3	$\dot{x}_3 = \begin{cases} \delta(t_f)\Delta x_3, & t = t_f \\ 0, & \text{otherwise} \end{cases}$	Δx_3
A_i^-	Impeller wear	A_i	$\dot{A}_i(t) = -w_{A_i} q_i(t)^2$	w_{A_i}
r_t^+	Thrust bearing wear	r_t	$\dot{r}_t(t) = w_t r_t \omega^2$	w_t
r_r^+	Radial bearing wear	r_r	$\dot{r}_r(t) = w_r r_r \omega^2$	w_r

Table 2: Fault signatures.

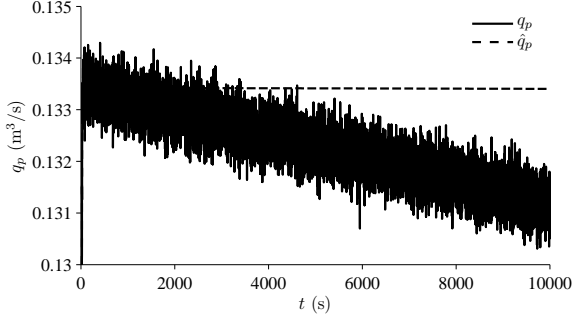
Faults	h_1	h_2	p_s	p_d	ω_p	q_p	T_t	T_r	T_o
A_{I1}^+	0-	0-	0-	0+	0-	0-	0-	0-	0-
A_{I2}^+	0-	0-	0-	0-	0+	0+	0+	0+	0+
x_1^-	0+	0-	0-	0-	0-	0-	0-	0-	0-
x_2^-	0+	0-	0+	0+	0-	0-	0-	0-	0-
x_3^+	0-	0+	0-	0+	0-	0-	0-	0-	0-
A_i^-	0+	0-	0+	0-	0-	0-	0-	0-	0-
r_t^+	0+	0-	0+	0-	0-	0-	0+	0+	0+
r_r^+	0+	0-	0+	0-	0-	0-	0+	0+	0+

4.2 Demonstration of Approach

We now present a detailed integrated diagnosis and prognosis scenario to illustrate the approach. In this scenario, impeller wear begins at $t = 0$ s with wear rate $w_A = 3 \times 10^{-3}$. A fault is detected at $t = 1380$ s, via a decrease in the pump flow q_p , shown in Fig. 4. The initial candidate list is $\{A_i^-, r_t^+, r_r^+, A_{I1}^+, x_1^-, x_2^-, x_3^+\}$. The fault signatures are given in Table 2. At 3279 s, an increase in h_1 and a decrease in h_2 are detected, eliminating a leak in the storage tank and a stuck fault in V_3 . At 4568 s, an increase in p_d is detected, which eliminates a stuck fault of V_2 , and at 7117 s, an increase in p_s is detected, eliminating a stuck fault of V_1 .

For our experiments, we had set adopted a policy to initiate fault identification once the number of fault candidates reduces to three or less, or if the qualitative fault isolation module has executed for 3000 s. We used particle filters with $N = 50$ particles for the nominal observer used for fault detection, and each faulty observer used for fault identification. The results of this experimental run is summarized in Table 3

Fig. 5 shows the (filtered) summed output errors for two fault candidates, impeller wear and thrust bearing wear. It is clear that by 7000 s impeller wear is the true candidate. The remaining fault candidates at this time have similar error to the error for thrust bearing wear, and can be eliminated, so the true fault is identified. Fig. 6 shows the estimated wear parameter estimate for impeller wear. Because the fault progression is so slow, by the end of the first fueling (at 10,000 s) the estimate is still converging. In further fuelings the estimate has converged with a small spread and remains fairly steady, due to a variance control algorithm pre-

Figure 4: Measured and predicted q_p values.

sented in (Daigle and Goebel, 2011) that dynamically modifies the random walk variance of the prediction algorithm to maintain a user-specified relative spread of the unknown fault parameters. The corresponding RUL predictions, made at the halfway point and the end of each fueling are shown in Fig. 7. By the third prediction point, the algorithm has converged and predictions remain within the desired accuracy window of 10%. The predictions were made assuming known future system inputs, so the uncertainty in the predictions is due solely to that resulting from the identification stage.

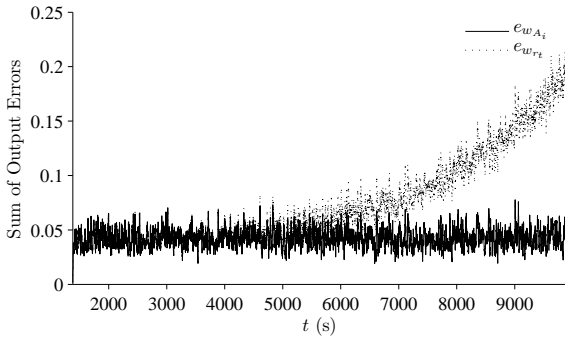
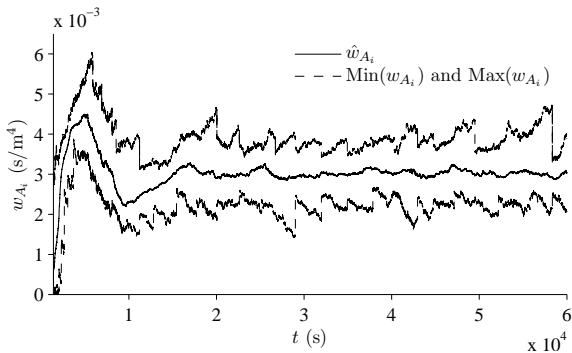


Figure 5: Sum of output errors for identification of impeller wear and thrust bearing wear fault candidates.

Figure 6: Estimated w_{A_i} values.

4.3 Simulation Results

Table 3 summarizes the results of several simulation experiments. The columns of the table represent the true fault; true injected value of the fault parameter ϕ_f ; Δk_d , the time in seconds to detect the fault; Δk_i , the time in seconds to isolate the true fault; the set of fault candidates after qualitative fault isolation; the estimated value of the fault parameter, ϕ_f , determined by the fault identifier; and \overline{RA} , the average relative accuracy (RA) over every 5000 s, where RA is defined as

$$RA_{k_P} = 100 \left(1 - \frac{|\text{RUL}_{k_P}^* - \widehat{\text{RUL}}_{k_P}|}{\widehat{\text{RUL}}_{k_P}} \right),$$

such that $\text{RUL}_{k_P}^*$ is the true RUL at time k_P , and $\widehat{\text{RUL}}_{k_P}$ is the mean of the prediction. For the abrupt faults, EOL is reached as soon as the fault is detected, and hence, \overline{RA} is not applicable. For the pump wear faults, however, the EOL is reached when certain thresholds are reached. The results indicate that fault detection and isolation times are fairly slow, due to slow progression of faults effects. However, we predict EOL with high \overline{RA} , usually ranging above 95%.

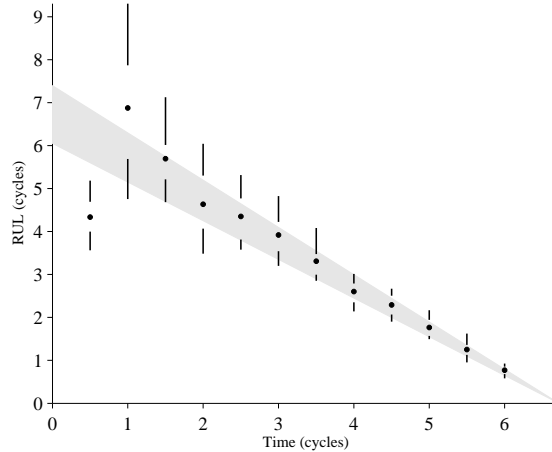


Figure 7: Predicted RUL of pump in the number of loading cycles (1 cycle = 10000 s). The mean is indicated with a dot and confidence intervals for 5% and 95% by lines. The gray cone depicts an accuracy requirement of 10%.

5 CONCLUSIONS

This paper presented an integrated diagnostic and prognostic framework. Our approach makes use of a common modeling paradigm to model both the nominal, as well as the fault progression models. We demonstrated our approach on a representative fuel transfer system. Our experimental results showed that our integrated diagnostic and prognostic approach diagnosed faults and prognosed the EOL accurately.

In future work, we would like to apply this approach to larger systems, to study the scalability of our diagnosis and prognosis scheme. Further, we would like to

Table 3: Diagnosis Results

True Fault	True ϕ_f	Δk_d	Δk_i	Fault Candidates	Estimated ϕ_f	RA
Nominal	N/A	∞	∞	\emptyset	N/A	N/A
A_{I1}^+	1.00×10^{-3}	77	77	$\Delta A_{I1} = 1.00 \times 10^3, e = 2.36 \times 10^{-3}$ $\Delta A_{I2} = 4.89 \times 10^{-4}, e = 4.01 \times 10^1$ $\Delta x_3 = 1.46, e = 1.08 \times 10^1$	$\Delta A_{I1} = 1.00 \times 10^{-3}$	N/A
A_{I2}^+	1.00×10^{-3}	236	236	$w_{A_i} = 1.08 \times 10^{-3}, e = 7.26$ $w_t = 2.82 \times 10^{-11}, e = 9.06$ $w_r = 2.92 \times 10^{-11}, e = 9.03$ $\Delta A_{I1} = 2.53 \times 10^{-4}, e = 9.06$ $\Delta A_{I2} = 1.01 \times 10^{-3}, e = 2.29 \times 10^{-3}$ $\Delta x_1 = 1.80, e = 9.02$ $\Delta x_2 = 2.08, e = 9.02$	$\Delta A_{I2} = 1.01 \times 10^{-3}$	N/A
x_1^-	-5.00×10^{-1}	0	1290	$\Delta x_1 = -5.00 \times 10^{-1}, e = 2.28 \times 10^{-3}$	$\Delta x_1 = -5.00 \times 10^{-1}$	N/A
x_2^-	-5.00×10^{-1}	0	1204	$\Delta x_2 = -5.00 \times 10^{-1}, e = 2.31 \times 10^{-3}$	$\Delta x_2 = -5.00 \times 10^{-1}$	N/A
x_3^+	5.00×10^{-1}	105	105	$\Delta x_3 = 5.09 \times 10^{-1}, e = 2.30 \times 10^{-3}$	$\Delta x_3 = 5.09 \times 10^{-1}$	N/A
A_i^-	3.00×10^{-3}	1379	7116	$w_{A_i} = 3.00 \times 10^{-3}, e = 2.50 \times 10^{-3}$ $w_t = 9.33 \times 10^{-13}, e = 1.63 \times 10^1$ $w_r = 2.27 \times 10^{-13}, e = 1.63 \times 10^1$	$w_{A_i} = 3.00 \times 10^{-3}$	96.19
r_t^+	8.00×10^{-11}	614	614	$w_t = 7.88 \times 10^{-11}, e = 2.63 \times 10^{-3}$ $w_r = 3.96 \times 10^{-12}, e = 1.58 \times 10^5$ $\Delta A_{I2} = 1.53 \times 10^{-6}, e = 1.85 \times 10^5$	$w_t = 7.88 \times 10^{-11}$	96.75
r_r^+	9.00×10^{-11}	748	748	$w_t = 7.12 \times 10^{-13}, e = 6.66 \times 10^4$ $w_r = 8.71 \times 10^{-11}, e = 2.70 \times 10^{-3}$ $\Delta A_{I2} = 1.47 \times 10^{-7}, e = 7.16 \times 10^4$	$w_r = 8.71 \times 10^{-11}$	91.28

expand the capability of this approach to hybrid systems. We would also like to enhance this approach to include multiple fault diagnosis and prognosis. Finally, we would like to investigate system level diagnosis and prognosis schemes.

REFERENCES

- (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.
- (Biswas and Mahadevan, 2007) G. Biswas and S. Mahadevan. A hierarchical model-based approach to systems health management. In *2007 IEEE Aerospace Conference*, March 2007.
- (Daigle and Goebel, 2010) M. Daigle and K. Goebel. Model-based prognostics under limited sensing. In *2010 IEEE Aerospace Conference*, March 2010.
- (Daigle and Goebel, 2011) M. Daigle and K. Goebel. Multiple damage progression paths in model-based prognostics. In *Proceedings of the 2011 IEEE Aerospace Conference*, March 2011.
- (Daigle *et al.*, 2010) M. J. Daigle, I. Roychoudhury, G. Biswas, and X. Koutsoukos. A comprehensive diagnosis methodology for complex hybrid systems: A case study on spacecraft power distribution systems. *IEEE Transactions on System, Man, and Cybernetics, Part A: Special issue on "Model-based Diagnosis: Facing Challenges in Real-world Applications"*, 4(5):917–931, September 2010.
- (Goodrich *et al.*, 2009) C. Goodrich, S. Narasimhan, M. Daigle, W. Hatfield, R. Johnson, and B. Brown. Applying model-based diagnosis to a rapid propellant loading system. In *Proceedings of the 20th International Workshop on Principles of Diagnosis*, pages 147–154, June 2009.
- (Hutchings, 1992) I. M. Hutchings. *Tribology: friction and wear of engineering materials*. CRC Press, 1992.
- (Kallesøe, 2005) C.S. Kallesøe. *Fault detection and isolation in centrifugal pumps*. PhD thesis, Aalborg University, 2005.
- (Lyshevski, 1999) S. E. Lyshevski. *Electromechanical Systems, Electric Machines, and Applied Mechatronics*. CRC, 1999.
- (Mosterman and Biswas, 1999) P. J. Mosterman and G. Biswas. Diagnosis of continuous valued systems in transient operating regions. *Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, 29(6):554–565, November 1999.
- (Roychoudhury, 2009) I. Roychoudhury. *Distributed Diagnosis of Continuous Systems: Global Diagnosis Through Local Analysis*. PhD thesis, Vanderbilt University, 2009.
- (Tu *et al.*, 2007) Fang Tu, S. Ghoshal, Jianhui Luo, G. Biswas, S. Mahadevan, L. Jaw, and K. Navarra. PHM integration with maintenance and inventory management systems. In *Proc. of the 2007 IEEE Aerospace Conference*, March 2007.
- (Wolfram *et al.*, 2001) A. Wolfram, D. Fussel, T. Brune, and R. Isermann. Component-based multi-model approach for fault detection and diagnosis of a centrifugal pump. In *Proceedings of the 2001 American Control Conference*, volume 6, pages 4443–4448, 2001.