# A behavioral approach to direct data-driven fault detection

Ivan Markovsky[a,b,*,1], Alba Muixí[b,c], Sergio Zlotnik[b,c] and Pedro Diez[b,c]

[a]*Catalan Institution for Research and Advanced Studies, Pg. Lluis Companys 23, Barcelona 08010, Spain*

[b]*Centre Internacional de Mètodes Numèrics en Enginyeria, Gran Capità, Edifici C1, Barcelona 08034, Spain*

[c]*Universitat Politècnica de Catalunya, LaCaN, DECA, Gran Capità, Edifici C1, Barcelona 08034, Spain*

**ARTICLE INFO**

**Abstract**

Efficient and reliable fault detection methods are needed for monitoring and evaluation of processes, *e.g.*, in structural health assessment. Existing methods rely on an a priori given model. Obtaining a model is nontrivial and is often the bottleneck in applications. Direct data-driven methods were recently developed in signal processing and control. These methods avoid the model identification step and may outperform state-of-the-art model-based methods. This paper presents a direct data-driven method for fault detection. The monitored process is modeled as a linear time-invariant system with an unobserved deterministic disturbance. The behavioral approach to systems theory is used in order to define a representation invariant measure for the distance between data and model. The main contribution of the paper is a method for computing the distance directly from offline and online data without parametric model identification. Another contribution is a direct data-driven input estimation. The method is validated empirically on simulated data of a flexible beam with a crack, where the objective is to detect the crack.

## 1. Introduction

Fault detection is a real-time monitoring problem aimed to determine based on data of a dynamical system if the system is in a "healthy" mode of operation, see Kopsaftopoulos and Fassois (2010); Chen, Frank, Kinnaert, Lunze and Patton (2001); Antoni, Kestel, Peeters, Leclère, Girardin, Ooijevaar and Helsen (2024). The healthy mode of operation, as well as possible "faulty" modes of operation, are specified by parametric models or data collected offline when the system is in the corresponding mode.

There are two fundamentally different approaches for fault detection. The first one is based on parametric model identification Ljung (1999). The existence of a fault is detected by monitoring the estimated model parameters. The rationale for this approach is that a fault causes a change in the model parameters. The second approach is based on a signal, call a *residual*, that measures the discrepancy between the data and the model. The occurrence of a fault is detected by thresholding the residual signal. The rationale for this approach is that the residual remains "small" (ideally zero) in the healthy mode of operation, and becomes "large" in the presence of a fault. Residuals are defined in terms of system representations. For linear time-invariant systems, the representation could be an impulse response, a transfer function, or a state-space. Possible residuals are the prediction error, output error, or equation error.

The approach in this paper is reminiscent of the residual thresholding approach. However, it uses a data–model discrepancy measure that is not based on a representation of the system and can be computed directly from data. The new *representation invariant* distance measure, is defined in the behavioral setting, where dynamical systems are viewed as sets of trajectories Willems (2007). This view is naturally suited for data-driven analysis and design. In contemporary machine learning language, the behavioral setting is non-parametric and unsupervised since the data does not have to be labeled into inputs and outputs.

---

✉ imarkovsky@cimne.upc.edu (I. Markovsky); alba.muixi@upc.edu (A. Muixí); sergio.zlotnik@upc.edu (S. Zlotnik); pedro.diez@upc.edu (P. Diez)

🌐 https://imarkovs.github.io (I. Markovsky); https://www.lacan.upc.edu/users/zlotnik/cv-sz-html/index.html (S. Zlotnik); https://www.lacan.upc.edu/author/pedro-diez/ (P. Diez)

ORCID(s): 0000-0001-0000-0000 (I. Markovsky); 0000-0002-4420-3366 (A. Muixí); 0000-0001-9674-8950 (S. Zlotnik); 0000-0001-6464-6407 (P. Diez)

---

In the behavioral setting, a natural choice for the discrepancy or *misfit* between the system and a signal is the orthogonal distance from the signal to the system, *i.e.*, the projection of the signal on the system. This operation is equivalent to errors-in-variables Kalman smoothing Markovsky and De Moor (2005). It is statistically optimal under the assumption that the signal is generated in the errors-in-variables setting, *i.e.*, the data is a trajectory of the system corrupted by measurement errors, which are assumed to be zero mean, white, Gaussian. The misfit attributes the lack of fit between the data and the model to the measurement errors. An alternative data–model discrepancy measure, called *latency*, attributes the lack of fit between the data and the model to an unobserved latent input signal, called *disturbance* Lemmerling and De Moor (2001). In the latency setting, the model describes the joint dynamics of the disturbance and the observed variables. The latency is defined as the size of the smallest disturbance that makes the observations compatible with the model. The latency computation corresponds to Kalman smoothing Willems (2004). The latency is statistically optimal in the auto-regressive moving-average exogenous (ARMAX) setting under suitable stochastic assumptions about the disturbance signal.

We define a new distance measure that combines misfit and latency with disturbance signal that is an unknown deterministic signal. Relaxing the assumptions about the disturbance allows us to apply the new distance in applications where there is no prior information about the disturbance. The new distance measure coincides with the misfit when the system has no disturbance. Also, the new distance measure is zero when there is no measurement error, *i.e.*, the signal is exact. As a byproduct of its computation, the disturbance signal is estimated. Thus, an independent contribution is a new input estimation method, see Gillijns and De Moor (2007); Eftekhar Azam, Chatzi and Papadimitriou (2015); Abooshahab, Alyaseen, Bitmead and Hovd (2022).

Based on the representation invariant distance measure, we propose a fault detection method. In addition to being representation invariant, the new distance measure and resulting fault detection method allow for direct computation from offline data of the system, bypassing the parametric model identification required by model-based methods. The approach of using data as a representation of the system is called *direct data-driven* and is successfully used in signal processing and control, where it is shown to have advantages over alternative model-based methods Markovsky (2015); Markovsky and Dörfler (2021); Markovsky, Huang and Dörfler (2023a).

The direct data-driven method proposed in the paper applies to data obtained from a transient response, forced response due to observed excitation signal, as well as forced response due to unobserved excitation signal. It can be computed efficiently in real-time and is validated on simulated data of a lumped mechanical system consisting of three masses connected by strings and dampers. As a realistic validation, we show the performance of the method for detection of a crack in a vibrating flexible beam.

The contributions of the paper are:

1. new measure for the distance between data and model that unifies the concepts of misfit and latency,
2. direct data-driven method for input estimation, and
3. end-to-end fault detection method.

Section 2 introduces the terminology, notation, and results from the behavioral systems theory that are used in the paper. The concepts of misfit, latency, and the new distance measure that unifies them are defined in Section 3. Section 4 presents the direct data-driven fault detection method based on the new distance measure. The method is validated in Section 5 on data of a lumped mechanical system as well as a distributed system—flexible beam with a crack. Details about the simulation of the beam with a crack are given in Appendix A.

## 2. Preliminaries and notation

We use the behavioral approach Polderman and Willems (1998); Willems (2007); Markovsky et al. (2023a). A real-valued $q$-variate signal $w$ with time axis $\mathcal{T} \subset \mathbb{R}$ is a map from $\mathcal{T}$ to $\mathbb{R}^q$. The set of signals $w : \mathcal{T} \to \mathbb{R}^q$ with $q$ variables is denoted by $(\mathbb{R}^q)^{\mathcal{T}}$. In this paper, the signals are discrete-time and $\mathcal{T} = \mathbb{N}$ — the time axis is the set of natural numbers. The *unit shift operator* is

$$(\sigma w)(t) := w(t + 1).$$

In the behavioral setting, a *dynamical system* $\mathcal{B}$ with $q$ variables is defined as a subset of the set of signals $(\mathbb{R}^q)^{\mathcal{T}}$. A system $\mathcal{B}$ is *linear* if $\mathcal{B}$ is a linear subspace of $(\mathbb{R}^q)^{\mathcal{T}}$ and *time-invariant* if $\mathcal{B}$ is invariant to the action of the shift operator, *i.e.*, $\sigma \mathcal{B} = \mathcal{B}$. The set of *linear time-invariant systems* with $q$ variables is denoted by $\mathcal{L}^q$.

The *restriction* of a signal $w \in (\mathbb{R}^q)^{\mathbb{N}}$ and a system $\mathcal{B} \subset (\mathbb{R}^q)^{\mathbb{N}}$ to the interval $1, \dots, T$ is denoted by $w|_T$ and $\mathcal{B}|_T$, respectively. The *restricted behavior* $\mathcal{B}|_T$ is a subspace of the set $(\mathbb{R}^q)^T$. When $\mathcal{B}$ is linear time-invariant,

$$\dim \mathcal{B}|_T = mT + n, \qquad \text{for all } T \geq \ell,$$

where $m$, $\ell$, and $n$ are natural numbers that are properties of the system $\mathcal{B}$:

- $m$ is the *number of inputs* (in an input/output representation of the system),

- $\ell$, called the *lag* of $\mathcal{B}$, is the minimal degree of a difference equation representation of $\mathcal{B}$, and

- $n$, called the *order* of $\mathcal{B}$, is the minimal total degree of a difference equation representation of $\mathcal{B}$.

The triple $(m, \ell, n)$ characterizes the *complexity* of $\mathcal{B} \in \mathcal{L}^q$. The set of *linear time-invariant systems* with $q$ variables and complexity bounded by $(m, \ell, n)$ is denoted by $\mathcal{L}^q_{(m,\ell,n)}$.

The variables $w$ can be partitioned into *inputs* $u$ (free variables) and *outputs* $y$ (dependent variables) via a permutation matrix $\Pi \in \mathbb{R}^{q \times q}$, *i.e.*, $w = \Pi \left[ \begin{smallmatrix} u \\ y \end{smallmatrix} \right]$. The inputs $u$ can be chosen freely while the outputs $y$ are uniquely defined by the model, the given inputs $u$, and the initial conditions. As shown in (Markovsky and Rapisarda, 2008, Lemma 1), the *initial conditions* for a trajectory $w = \bigl(w(1), \dots, w(T)\bigr)$ can be specified by $T_{\text{ini}} \geq \ell$ "past" samples

$$w_{\text{ini}} = \bigl(w(-T_{\text{ini}} + 1), \dots, w(0)\bigr).$$

A partitioning of the variables into inputs and outputs is not unique. In the context of fault detection, we use an input/output partitioning in order to model user defined excitation signals and disturbances.

The restricted behavior $\mathcal{B}|_T$ of a linear time-invariant system $\mathcal{B} \in \mathcal{L}^q$ is an

$$r := \dim \mathcal{B}|_T = Tm + n$$

dimensional subspace and, therefore, it can be represented by a basis

$$\mathcal{B}|_T = \text{image } B, \quad \text{where } B := \begin{bmatrix} b^1 & \cdots & b^r \end{bmatrix} \in \mathbb{R}^{qT \times r}. \tag{B}$$

With some abuse of the terminology, we refer to the matrix $B$ of the basis vectors as the basis. The representation (B) of $\mathcal{B}|_T$ is *nonparameteric*. A trajectory $w \in \mathcal{B}|_T$ is specified using the data-driven representation (B) via the equation $w = Bg$, where $g \in \mathbb{R}^r$. For $T \geq \ell(\mathcal{B}) + 1$, the basis $B$ for the finite-horizon behavior $\mathcal{B}|_T$ uniquely defines the system $\mathcal{B}$, see (Markovsky and Dörfler, 2023, Lemma 13).

Consider a finite trajectory $w_{\text{d}} \in \mathcal{B}|_{T_{\text{d}}}$ (the subscript index "d" stands for "data") of a bounded complexity linear time-invariant system $\mathcal{B} \in \mathcal{L}^q_{(m,\ell,n)}$. Theorem 17 from Markovsky and Dörfler (2023) gives conditions under which a basis $B$ for $\mathcal{B}|_T$ can be obtained from the data $w_{\text{d}}$. Define the Hankel matrix $\mathcal{H}_T(w_{\text{d}})$ with depth $T$

$$\mathcal{H}_T(w_{\text{d}}) := \begin{bmatrix} w|_T & (\sigma w)|_T & \cdots & (\sigma^{T_{\text{d}}-T} w)|_T \end{bmatrix} \in \mathbb{R}^{qT \times (T_{\text{d}}-T+1)}.$$

The result of (Markovsky and Dörfler, 2023, Theorem 17) states that, for any $T \geq \ell$, the finite horizon behavior $\mathcal{B}|_T$ of the data-generating system is equal to the image of the Hankel matrix $\mathcal{H}_T(w_{\text{d}})$ of the data,

$$\mathcal{B}|_T = \text{image } \mathcal{H}_T(w_{\text{d}})$$

if and only if

$$\text{rank } \mathcal{H}_T(w_{\text{d}}) = mT + n. \tag{GPE}$$

The condition (GPE), called *generalized persistency of excitation*, is verifiable from the data $w_{\text{d}}$ and the model's complexity $(m, \ell, n)$. This result is a generalization of the *fundamental lemma* Willems, Rapisarda, Markovsky and De Moor (2005). For detailed discussion on the similarities and differences between the generalized persistency of excitation and the conditions of the fundamental lemma, see Markovsky, Prieto-Araujo and Dörfler (2023b).

Linear time-invariant systems can be represented in different ways by equations. The most popular ones—convolution, transfer function, and state-space—assume a priori given input/output partitioning of the variables. In this section, we do not review parametric representations because they are not used in the paper. For more details on the behavioral approach to systems theory, it's relation to the classical input/output approach, and its relevance to direct data-driven signal processing and control, we refer the reader to Markovsky et al. (2023a).

## 3. Data–model discrepancy measures

Stochastic system identification problems and methods can be classified into errors-in-variables Söderström (2018) and auto-regressive moving-average exogenous (ARMAX). Deterministic counterparts of the likelihood functions in the errors-in-variables and the ARMAX settings are, respectively, the misfit and the latency. In this section, we propose a new distance measure that combines and generalizes the misfit and the latency.

We model the disturbance as an unknown deterministic input, *i.e.*, without imposing any prior assumptions about it. Under conditions on the system however the disturbance can be inferred from the observed variables. In this section, we assume that the system is given. In the next section, we show how the new distance measure can be estimated directly from data.

The misfit between a signal $w \in (\mathbb{R}^q)^T$ and a system $\mathcal{B} \subset \mathcal{L}^q$ is defined as the minimum norm perturbation of $w$ that makes the perturbed signal $\widehat{w}$ consistent with the system $\mathcal{B}$, *i.e.*,

$$\text{misfit}(w, \mathcal{B}) := \min_{\widehat{w} \in \mathcal{B}|_T} \|w - \widehat{w}\|. \tag{M}$$

The misfit$(w, \mathcal{B})$ is the likelihood of $w$ given $\mathcal{B}$ in the *errors-in-variables setting:*

$$w = \overline{w} + \widetilde{w}, \qquad \text{where } \overline{w} \in \mathcal{B}|_T \text{ and } \widetilde{w} \sim \text{N}(0, s^2 I_q), \tag{EIV}$$

*i.e.*, $\overline{w}$ is an exact trajectory of $\mathcal{B}$ and $\widetilde{w}$ is a stochastic process with zero-mean, white, Gaussian distribution with a covariance matrix $s^2 I_q$, where $I_q$ is the identity matrix of size $q$.

**Lemma 1.** *Assuming that the data $w$ is generated in the errors-in-variables setting, misfit$(w, \mathcal{B}) \leq \|\widetilde{w}\|$.*

*Proof.* For $w = \overline{w} + \widetilde{w}$, where $\overline{w} \in \mathcal{B}|_T$, we have that misfit$(w, \mathcal{B}) = \min_{\Delta w \in \mathcal{B}|_T} \|\widetilde{w} - \Delta w\|$. Since $\Delta w = 0 \in \mathcal{B}|_T$, misfit$(w, \mathcal{B}) \leq \|\widetilde{w}\|$. $\square$

An alternative way of measuring the discrepancy between data and model is an unobserved signal $e \in (\mathbb{R}^{n_e})^{\mathbb{N}}$ acting on the system. In this case, the system $\mathcal{B}$ describes the extended signal $(e, w) \in (\mathbb{R}^{n_e+q})^{\mathbb{N}}$. The *latency* of $w \in (\mathbb{R}^q)^T$, given $\mathcal{B} \in \mathcal{L}^{n_e+q}$ is defined as

$$\text{latency}(w, \mathcal{B}) = \min_{(\widehat{e}, w) \in \mathcal{B}|_T} \|\widehat{e}\|. \tag{L}$$

For latency$(w, \mathcal{B})$ to be well-defined, problem (L) should have a unique solution. A necessary and sufficient condition for existence and uniqueness of solution is that $w$ is compatible with $\mathcal{B}$, *i.e.*, $w \in \Pi_w \mathcal{B}|_T$, where $\Pi_w$ is the projection of $(e, w)$ on the $w$ component. The condition is satisfied when 1) $e$ is an input of $\mathcal{B}$ and 2) $n_e = p$, where $p$ is the number of outputs of $\mathcal{B}$. These assumptions are standard in the ARMAX setting Ljung (1999):

$$(\overline{e}, w) \in \mathcal{B}|_T, \tag{ARMAX}$$

where the disturbance $\overline{e}$ is a zero-mean, white, Gaussian process. The latency is the likelihood of $w$ given $\mathcal{B}$ in the ARMAX setting. The following statement follows directly from the assumption that $(\overline{e}, w) \in \mathcal{B}|_T$, and the definition of latency$(w, \mathcal{B})$.

**Lemma 2.** *Assuming that the data $w$ is generated in the ARMAX setting, latency$(w, \mathcal{B}) \leq \|\overline{e}\|$.*

Consider next the missing data estimation problem: Given a system $\mathcal{B} \in \mathcal{L}^{n_e+q}$ with variables partitioned as $(e, w)$ and a signal $w \in (\mathbb{R}^q)^T$, find $e$, such that $(e, w) \in \mathcal{B}|_T$. The goal is to achieve *exact recovery* of $e$ from $w$, *i.e.*, there should be a unique signal $e \in (\mathbb{R}^{n_e})^T$ that is compatible with the data $w$ and the model $\mathcal{B}$. The following result from Markovsky and Dörfler (2022) gives necessary and sufficient conditions for exact recovery.

**Lemma 3.** *Consider a system $\mathcal{B} \in \mathcal{L}^{n_e+q}_{(m,\ell,n)}$, let $w \in \Pi_w \mathcal{B}|_T$, and let $B$ be a basis for $\mathcal{B}|_T$. There is a unique $e \in (\mathbb{R}^{n_e})^T$, such that $(e, w) \in \mathcal{B}|_T$ if and only if*

$$\text{rank } \Pi_w B = mT + n. \tag{A}$$

A necessary condition for unique recovery is $p > n_e$, where $p$ denotes the number of outputs of $\mathcal{B}$. Note that in the ARMAX setting unique recovery of $e$ is not possible.

In the missing data estimation problem, the data $w$ may be corrupted by measurement noise. Then, generically, there is no $e \in (\mathbb{R}^{n_e})^T$, such that $(e, w) \in \mathcal{B}|_T$. In this case, under the assumptions of Lemma 3, we choose the signal $\hat{e}$ that achieves the best in the least-squares sense fit to the data $w$:

$$\text{dist}(w, \mathcal{B}) := \min_{(\hat{e}, \hat{w}) \in \mathcal{B}|_T} \|w - \hat{w}\|. \tag{dist}$$

Problem (dist) is a generalization of (M). Indeed, under the assumptions of Lemma 3, when there is no latent input, (dist) coincides with (M). Moreover, dist$(w, \mathcal{B})$ is the likelihood of $w$, given $\mathcal{B}$, when

$$w = \overline{w} + \widetilde{w}, \quad \text{where } (\overline{e}, \overline{w}) \in \mathcal{B}|_T \text{ and } \widetilde{w} \sim \text{N}(0, s^2 I_q), \tag{EIV-ARMAX}$$

for some $\overline{e} \in (\mathbb{R}^{n_e})^T$ and measurement noise $\widetilde{w}$ that is a zero-mean, white, Gaussian process.

**Lemma 4.** *Under the assumptions of Lemma 3, if the data $w$ is generated in the (EIV-ARMAX) setting,* dist$(w, \mathcal{B}) \leq \|\widetilde{w}\|$.

*Proof.* By Lemma 3, we have that

$$\text{dist}(w, \mathcal{B}) = \text{misfit}(w, \Pi_w \mathcal{B}). \tag{dist $\leftrightarrow$ misfit}$$

Lemma 4 follows then from Lemma 1. □

The basic idea of the direct data-driven fault detection method proposed in the paper is to compute and monitor in real-time the distance dist$(w, \mathcal{B})$ between the $T$ most recent measurements $w$ obtained from the process and a predefined horizon-$T$ behavior $\mathcal{B}$. The horizon $T$ in the computation of the distance measure is a hyper-parameter of the fault detection method. In theory, for exact data generated by a bounded complexity linear time-invariant system, it is sufficient to choose $T$ larger than the lag $\ell$ of the system. In the presence of noise, however, $T$ should be chosen as larger as possible in order to achieve noise smoothing. On the other hand, large $T$ compromises the speed of the fault detection method. Thus, the choice of $T$ involves a speed vs accuracy trade-off.

Next, we outline the proposed end-to-end fault detection method based on the computation of dist$(w, \mathcal{B})$.

## 4. Direct data-driven fault detection

In the data-driven fault detection problem considered, the monitored process is a bounded complexity linear time-invariant system with variables $(e, w)$. Its nominal behavior $\mathcal{B}^0$ is implicitly specified by offline data $(e_d^0, w_d{}^0) \in \mathcal{B}^0|_{T_d}$, that satisfies the generalized persistency of excitation condition (GPE). Possible faulty behaviors $\mathcal{B}^1, \ldots, \mathcal{B}^N$ are also specified by offline data $(e_d^i, w_d{}^i) \in \mathcal{B}^i|_{T_d}$ for $i = 1, \ldots, N$ that satisfy the generalized persistency of excitation condition. The fault detection problem aims to check if an observed signal $w \in (\mathbb{R}^q)^T$, generated in the (EIV-ARMAX) setting, is compatible with $\mathcal{B}^0$. If it isn't, then a possible fault among the predefined options $\mathcal{B}^1, \ldots, \mathcal{B}^N$ is localized. This is done by computing and comparing the distance measures $d_i := \text{dist}(w, \mathcal{B}^i)$ from the data $w$ to the nominal behavior $\mathcal{B}^0$ and the possible faulty behaviors $\mathcal{B}^1, \ldots, \mathcal{B}^N$.

The method proposed has two phases:

1. using the offline data $(e_d^i, w_d{}^i)$ and a complexity specification $(m, \ell, n)$, find orthonormal bases $B^i$ for $\mathcal{B}^i|_T$, and
2. using the online data $w \in (\mathbb{R}^q)^T$, compute $d_i := \text{dist}(w, \mathcal{B}^i)$, for $i = 0, 1, \ldots, N$.

In phase 1, the bases $B^i$ can be computed by low-rank approximation, *i.e.*, truncation of the singular value decomposition of the Hankel matrices $\mathcal{H}_T(w_d^i)$ to the theoretical rank $mT + n$. Although this method is computationally inexpensive and easy to implement, it does not preserve the shift-invariant structure and is statistically suboptimal. Alternatively, a Hankel structured low-rank approximation can be used, however, it leads to a nonconvex optimization problem, for details about this approach, see Markovsky (2019).

Problem (M) is a projection of $w$ on the subspace $\mathcal{B}|_T$. Thus, with $B$ being an orthonormal basis for $\mathcal{B}|_T$, we have

$$\text{misfit}(w, \mathcal{B}) = \min_g \|w - Bg\| = \sqrt{w^\top (I_{qT} - BB^\top) w}.$$

For the computation of the new measure $\text{dist}(w, \mathcal{B})$, note that by assumption (A), $B_w := \Pi_w B$ is a basis for $\Pi_w \mathcal{B}|_T$. However, it need not be an orthonormal basis. Then, using (dist ↔ misfit), we have that

$$\text{dist}(w, \mathcal{B}) = \min_g \|w - B_w g\| = \sqrt{w^\top (I_{qT} - B_w B_w^+) w},$$

where $B_w^+$ is the pseudo-inverse of $B_w$. Pre-orthogonalizing $B_w$ in the offline step, the cost for the evaluation of the distance $\text{dist}(w, \mathcal{B})$ in the online step is $O(T)$.

## 5. Empirical validation

The data-driven computation of the distance (dist) and estimation of the unobserved input signal $e$ are implemented in Matlab. The fault detection method based on (dist) is then empirically validated. First, we validate the method on data generated by a lumped system satisfying the assumption of bounded complexity linear time-invariant dynamics. In this example, the aim is to detect a fault, representing a change of a physical parameter of the system—mass of a body, spring constant, or damping coefficient. Then, we validate the method on data generated by a distributed system—a flexible beam. The system is defined by a partial differential equation in time and 2D space. The equation models a flexible beam with a crack. The aim in this case is to detect the presence of a crack from the observed vibrations at points on the surface of the beam. The distributed system does not satisfy the theoretical assumptions of the method for bounded complexity LTI dynamics.

### 5.1. Lumped system

In this section, we validate the method on data generated by a mechanical system consisting of three masses interconnected via springs and dampers, as shown in Figure 1. The system is excited by an external force $u$, applied on the first mass, and a disturbance $e$ acting on the third mass. Both $u$ and $e$ are generated as random independent and uniformly distributed in the interval $[0, 1]$ processes. The nominal system $\mathcal{B}^0$ has the following parameter values:

$$m_1 = m_2 = m_3 = 10, \ k_1 = k_2 = k_3 = 1, \ \text{and} \ b_1 = b_2 = b_3 = 0.5. \tag{PAR}$$

The observed signals are the force $u$ and the positions of the three masses. The hyper-parameters of the method—the horizon $T$ and the model's lag $\ell$—are chosen as $T = 100$ and $\ell = 2$. Note that the data-generating system's lag is 2, so that $\ell = 2$ is the theoretically optimal as well as computationally most economical choice. With the chosen discretization step, the transient process of the data-generating system is sufficiently decayed after 100 samples, so that $T = 100$ is a suitable choice for the time-horizon.
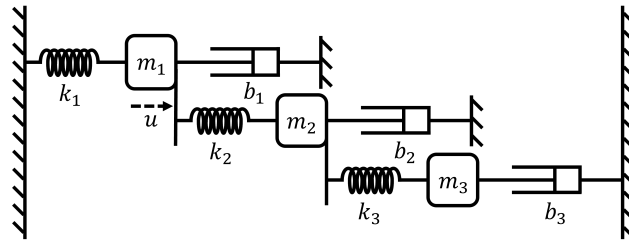


**Figure 1**: The empirical validation is done on an interconnected mass-spring-damper system with an external force $u$ applied on the first mass $m_1$ and a disturbance $e$ acting on the third mass $m_3$. The observed signals are the force $u$ and the positions of the three masses. The faults correspond to changes in the model parameters $m_1, k_1, b_1$ by 10% from their nominal value.

Three fault scenarios are considered. In each scenario, one of the parameter $m_1, k_1, b_1$ is increased by 10% of its nominal value. Data is collected from the nominal and the faulty systems in the following setups:

1. *transient process* — the input $u$ and the disturbance $e$ are set to zero and the output is generated from a random initial condition,
2. (EIV) — a random input $u$ signal acts on the system, while the disturbance $e$ is set to zero, and
3. (EIV-ARMAX) — both random observed input $u$ and unobserved disturbance $e$ act on the system.

In setup 1, the observed variables $w$ are the positions of the three masses. In setups 2 and 3, $w$ consists of the external force $u$ and the positions of the three masses. In all setups, the variables $w$ are measured with additive noise, which is simulated as a zero-mean, white, Gaussian process.

The distance measures $d_k := \text{dist}(w, \mathcal{B}^k)$ from observed data $w$ to the nominal and faulty behaviors are computed for increasing levels of the measurement noise variance $s$. The results shown in Figure 2 are for nominal data $w$, *i.e.*, data generated from the nominal model $\mathcal{B}^0$. The distances $d_k$ are computed for the three simulation setups—transient, EIV, and EIV-ARMAX—as a function of the noise level $s$ (the standard deviation of the additive measurement noise) averaged over 100 Monte-Carlo repetitions of the experiment. The distance $d_0$ from the nominal data to the nominal model serves as a reference for the distances $d_1, d_2, d_3$ from the nominal data to the faulty models $\mathcal{B}^1, \mathcal{B}^2, \mathcal{B}^3$. Without noise, *i.e.*, for $s = 0$, the distance $d_0 = 0$ and $d_i > 0$ for all $i = 1, 2, 3$. This indicates that the data comes from the nominal model. The larger the gap between $d_0$ and $d_i$ is, the easier and more reliable the fault detection is. With noise, distance $d_0 > 0$, however, $d_0$ remains the smallest among the evaluated distances. Thus, selecting the mode of operation related to the smallest distance, we can correctly conclude that the data is generated by the nominal model, *i.e.*, there is no fault. In the simulation examples, the gap between $d_0$ and $d_1$ is bigger than the one between $d_0$ and $d_3$. Thus, it is "easier" to detect faults related to a change of the spring coefficient $k_1$ than the damping coefficient $b_1$.
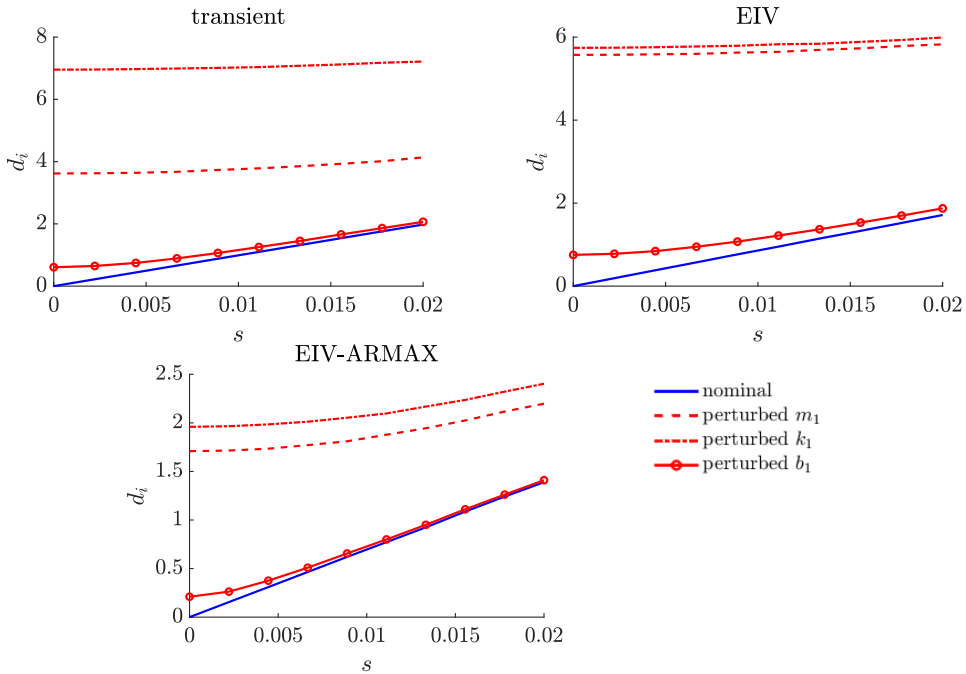


**Figure 2**: The smallest distances $d_k$ from the data $w$ to a behavior $\mathcal{B}^k$ identifies the current mode of operation. The results in the figure are for nominal data $w$, *i.e.*, data generated from $\mathcal{B}^0$, and the smallest distance is $d_0$ in all simulation setups—transient, EIV, and EIV-ARMAX—and for all noise level $s$. The margin between $d_0$ and $d_1, d_2, d_3$ indicates the robustness of the method—the larger the margin, the easier is to distinguish the modes.

Condition (A) of Theorem 3 is not satisfied for the system in Figure 1 (rank $\Pi_w \mathcal{B}^0 = 105$ while $m100 + n = 106$), so that exact recovery of the force $e$ from the observed positions $w$ of the masses is not possible. Nevertheless, the data-driven computation of the distance measure $\text{dist}(w, \mathcal{B})$ is well defined and fault detection based on it is possible.

## 5.2. Distributed system

Next, we validate the method on data obtained from a flexible beam, where the fault is a crack in the beam. First, we describe the simulation setup. Further details about the numerical simulation of the flexible beam with a crack are given in the appendix. Then, we present the results obtained with the data-driven fault detection method proposed in the paper. Zero-mean, white, Gaussian measurement noise is added to the simulated data in order to model measurement noise in a real-life experiment and test the robustness of the method.
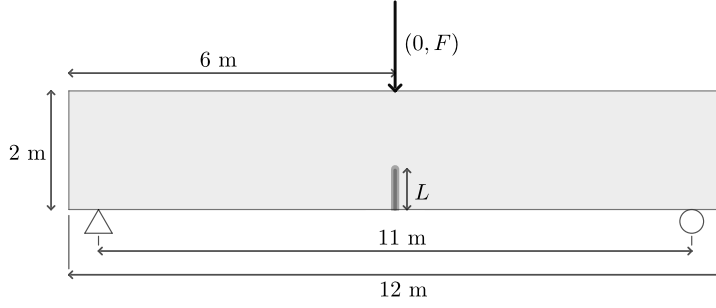
**Figure 3**: The simulation setup for the experiment with a distributed system is a flexible beam placed on two supports and subjected to a vertical point load $F$ at its top midpoint. A notch with length $L$ in the bottom middle of the beam represents a crack that we aim to detect from observed displacements.

**Table 1**
The fault scenarios correspond to different notch lengths $L$ and losses of stiffness.

| # | Length $L$ [m] | Loss of stiffness [%] |
|---|---|---|
| 0 | – | 0 |
| 1 | 0.7 | 100 |
| 2 | 0.7 | 36 |
| 3 | 0.2 | 100 |
| 4 | 0.2 | 36 |

We consider a two-dimensional model of a beam in a three-point bending setup. The beam is placed on two supports and subjected to a vertical point load $F$ at its top midpoint, as shown in Figure 3. The beam is pinned at the left support, while the right support acts as a roller that constrains only the vertical displacement. The objective is to detect structural damage from observation of the displacement at some points of the beam's surface.

The material is homogeneous, isotropic and linearly elastic with mass-proportional damping. The material parameters are the Young's modulus $E = 5$ GPa, the Poisson's ratio $\nu = 0.15$ and the mass density $\rho = 2450$ kg/m$^3$. The damping ratio is set to 2.5%. In the nominal case, the beam is undamaged. For the faults, we introduce a notch of length $L$ at the bottom midpoint of the beam. For this purpose, we use a phase-field variable $d$ that degrades the stiffness of the material along the notch, see Ambati, Gerasimov and De Lorenzis (2015).

The beam is simulated in the nominal and four fault scenarios, see Table 1. The fault scenarios account for a partial (36%) or complete (100%) loss of stiffness, representing a crack. We aim to detect the presence of a defect in the beam. Fault 1 is the most severe defect and case 4 the most subtle one.

For each experiment, the data consists of the displacements obtained from a finite element simulation of the model. The generalized $\alpha$-method temporal scheme is used for the simulation, Chung and Hulbert (1993). Numerical implementation details are given in Appendix A.

First, we perform a data collection experiment with a pulse excitation. The point load is increased linearly up to $F = -10^6$ N over 0.04 s, maintained at this value until $t = 0.16$ s, and then linearly decreased to zero over another 0.04 s. After being unloaded, the beam is left to vibrate freely. The simulation is run for a total time of 2 s using a uniform time discretization with 1000 steps. The data obtained after $t = 0.24$ s (*i.e.*, after the load is released) is refered to as *free vibration* data.

Second, we perform an experiment with a random excitation applied in the middle top of the beam. In this case, the external loading $F$ is piecewise constant with values sampled from a uniform distribution in the range $[-5, 5] \cdot 10^6$ N. The load value is kept constant for 1 s. The simulation spans a total time of 10 s, discretized into 2000 time steps.

The external force is measured in Newtons and varies in the interval $[-10^7, 10^7]$. The horizontal and vertical displacement at the grid points are measured in meters and vary in the interval $[-10^{-3}, 10^{-3}]$. In order to avoid numerical problems, we preprocess the data scaling all variables so that they are normalized in the interval $[-1, 1]$.

The number of measurement points and their optimal distribution on the structure for the fault detection task at hand is an interesting problem that we do not consider in this paper. As a practical constraint, we consider that positions and

accelerations can be measured only on the surface of the beam. Thus, for detection of a crack in the middle bottom of the beam, we select the surface points of the simulation grid that are immediately to the left and to the right of the middle bottom point. We conjecture that these points give us the most "informative data" for detecting the crack. The simulated data from the flexible beam is not consistent with the assumption (EIV-ARMAX) for the method. The inexactness of the data due to the infinite-dimensional nature of the distributed system can be attributed to "noise" however this noise is deterministic in nature. We refer to it as "modeling error," in order to distinguish it from the artificially added measurement error $\widetilde{w}$, which is a stochastic process. The hyper-parameters of the method—the distance computation horizon $T$ and the model's lag $\ell$—are fixed in all simulations as $T = 100$ and $\ell = 2$. As in Section 5.1, the results presented are the distances $d_0, d_1, \ldots, d_4$ between nominal data $w$ and the models $\mathcal{B}^0, \mathcal{B}^1, \ldots, \mathcal{B}^4$.

Both the offline and online free-response data has 4 variables—the horizontal and the vertical displacements at the selected points. The offline random excitation data has 5 variables—the excitation signal $u$ as well as the horizontal and vertical displacements at the selected points—while the online data has 4 variables—the horizontal and the vertical displacements at the selected points. For free response data, the distance measure (dist) is identical to the misfit (M). The results averaged over 100 Monte-Carlo repetitions are shown in Figure 4. Although noise-free data from the vibrating beam is not exact due to modeling errors, the distance measure $d_0$ is of the order of 1e-6. This suggests that the distributed system has dynamics that is very well approximated by a bounded complexity linear time-invariant system. Consequently, the results are qualitatively similar to the ones for the lumped system in Section 5.1.
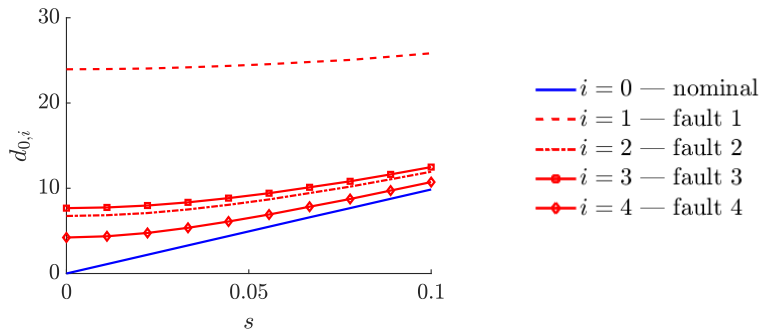


**Figure 4**: Distances $d_k$ from the nominal data $w$ to the corresponding behaviors $\mathcal{B}^k$ as a function of the noise level $s$ for the free vibration flexible beam. The results are qualitatively similar as the ones for a lumped linear time-invariant system.

For exact data obtained under an excitation signal $e$, the distance $d_0$ from exact nominal data to a linear time-invariant system $\widehat{\mathcal{B}}^0$ with lag $\ell = 2$ is again of the order of 1e-6. The plot in Figure 5 shows the estimate $\widehat{e}$ of $e$. The matrix $\Pi_w B$ in condition (A), however, is ill-conditioned. This leads to poor estimation of $e$ in case of noisy data. As in the case of a lumped system, the ill-conditioning of the estimation of $e$ does not affect the evaluation of the distance (dist). Consequently, the results are qualitatively similar to the ones for the lumped system.
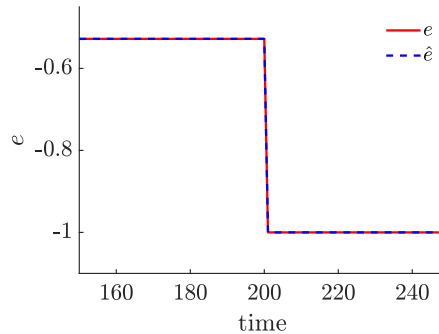


**Figure 5**: With exact data, the estimate $\widehat{e}$ computed by (dist) matches the excitation signal $e$ in all but the last sample.

## 6. Conclusions

We presented an approach for direct data-driven fault detection based on a new measure of distance for data to a system, where can be implicitly specified by data. The theoretical setting assumes bounded complexity linear time-invariant dynamics. The discrepancy between the data and the system is attributed to measurement errors as well as an unobserved disturbance signal. As a byproduct of computing the distance, the method estimates the disturbance signal—a result of independent interest. The method was empirically validated on both a lumped mechanical system and a distributed system for detection of a crack in a vibrating beam. Future work includes uncertainty quantification, input selection, and the application of the method in biomedical engineering for acute ischaemic stroke prediction.

## References

Abooshahab, M., Alyaseen, M., Bitmead, R., Hovd, M., 2022. Simultaneous input & state estimation, singular filtering and stability. Automatica 137, 110017. doi:10.1016/j.automatica.2021.110017.

Ambati, M., Gerasimov, T., De Lorenzis, L., 2015. A review on phase-field models of brittle fracture and a new fast hybrid formulation. Comput Mech 55, 383–405.

Amor, H., Marigo, J.J., Maurini, C., 2009. Regularized formulation of the variational brittle fracture with unilateral contact: numerical experiments. J Mech Phys Solids 57, 1209–1229.

Antoni, J., Kestel, K., Peeters, C., Leclère, Q., Girardin, F., Ooijevaar, T., Helsen, J., 2024. On the design of optimal health indicators for early fault detection and their statistical thresholds. Mechanical Systems and Signal Processing 218, 111518.

Borden, M.J., Verhoosel, C.V., Scott, M.A., Hughes, T.J.R., Landis, C.M., 2012. A phase-field description of dynamic brittle fracture. Comput Methods Appl Mech Eng 217-220, 77–95.

Chen, J., Frank, P.M., Kinnaert, M., Lunze, J., Patton, R.J., 2001. Fault Detection and Isolation. Springer London. chapter 9. pp. 191–207. doi:10.1007/978-1-4471-0349-3_9.

Chung, J., Hulbert, G.M., 1993. A time integration algorithm for structural dynamics with improved numerical dissipation: The generalized-$\alpha$ method. Journal of Applied Mechanics 60, 371–375. doi:10.1115/1.2900803.

Eftekhar Azam, S., Chatzi, E., Papadimitriou, C., 2015. A dual Kalman filter approach for state estimation via output-only acceleration measurements. Mechanical Systems and Signal Processing 60–61, 866–886. doi:https://doi.org/10.1016/j.ymssp.2015.02.001.

Gillijns, S., De Moor, B., 2007. Unbiased minimum-variance input and state estimation for linear discrete-time systems. Automatica 43, 111–116. doi:10.1016/j.automatica.2006.08.002.

Kopsaftopoulos, F., Fassois, S., 2010. Vibration based health monitoring for a lightweight truss structure: Experimental assessment of several statistical time series methods. Mechanical Systems and Signal Processing 24, 1977–1997. doi:https://doi.org/10.1016/j.ymssp.2010.05.013.

Langtangen, H.P., Logg, A., 2017. Solving PDEs in Python. The FEniCS Tutorial. Springer.

Lemmerling, P., De Moor, B., 2001. Misfit versus latency. Automatica 37, 2057–2067.

Ljung, L., 1999. System Identification: Theory for the User. Prentice-Hall.

Markovsky, I., 2015. Comparison of adaptive and model-free methods for dynamic measurement. IEEE Signal Proc. Lett. 22, 1094–1097. doi:10.1109/LSP.2014.2388369.

Markovsky, I., 2019. Low-Rank Approximation: Algorithms, Implementation, Applications. Springer. doi:10.1007/978-3-319-89620-5.

Markovsky, I., De Moor, B., 2005. Linear dynamic filtering with noisy input and output. Automatica 41, 167–171.

Markovsky, I., Dörfler, F., 2021. Behavioral systems theory in data-driven analysis, signal processing, and control. Annual Reviews in Control 52, 42–64. doi:10.1016/j.arcontrol.2021.09.005.

Markovsky, I., Dörfler, F., 2022. Data-driven dynamic interpolation and approximation. Automatica 135, 110008. doi:10.1016/j.automatica.2021.110008.

Markovsky, I., Dörfler, F., 2023. Identifiability in the behavioral setting. IEEE Trans. Automat. Contr. 68, 1667–1677. doi:10.1109/TAC.2022.3209954.

Markovsky, I., Huang, L., Dörfler, F., 2023a. Data-driven control based on behavioral approach: From theory to applications in power systems. IEEE Control Systems Magazine 43, 28–68. doi:10.1109/MCS.2023.3291638.

Markovsky, I., Prieto-Araujo, E., Dörfler, F., 2023b. On the persistency of excitation. Automatica , 110657doi:10.1016/j.automatica.2022.110657.

Markovsky, I., Rapisarda, P., 2008. Data-driven simulation and control. Int. J. Control 81, 1946–1959. doi:10.1080/00207170801942170.

Polderman, J., Willems, J.C., 1998. Introduction to Mathematical Systems Theory. Springer-Verlag.

Söderström, T., 2018. Errors-in-Variables Methods in System Identification. Springer.

Willems, J.C., 2004. Deterministic least squares filtering. J. Econometrics 118, 341–373.

Willems, J.C., 2007. The behavioral approach to open and interconnected systems: Modeling by tearing, zooming, and linking. IEEE Control Systems Magazine 27, 46–99.

Willems, J.C., Rapisarda, P., Markovsky, I., De Moor, B., 2005. A note on persistency of excitation. Systems & Control Lett. 54, 325–329. doi:10.1016/j.sysconle.2004.09.003.

## A. Appendix: Numerical simulation

This section presents the model equations and describes the numerical simulation for the benchmark in Section 5.2. In the domain $\Omega = [0, 12] \times [0, 2]$ m$^2$ occupied by the beam, the model for the displacement field $\boldsymbol{y}$ is

$$\rho \ddot{\boldsymbol{y}} + \gamma \dot{\boldsymbol{y}} - \boldsymbol{\nabla} \cdot \boldsymbol{\sigma}(\boldsymbol{y}) = \boldsymbol{F}, \quad \text{with } \boldsymbol{\sigma} = g(d)\, \boldsymbol{C} : \boldsymbol{\varepsilon}, \tag{1}$$

where $\rho$ is the mass density, $\boldsymbol{\sigma}(\boldsymbol{\varepsilon})$ is the stress tensor, $\boldsymbol{C}$ is a fourth-order tensor depending on the Lamé parameters and $\boldsymbol{\varepsilon}(\boldsymbol{y})$ is the small strain tensor. The damping coefficient $\gamma$ is taken as $\gamma = 2\rho\omega_0\zeta$, with damping ratio $\zeta = 2.5 \cdot 10^{-2}$ and $\omega_0$ the lowest natural angular frequency. The modal analysis of the beam leads to $\omega_0 = 119.13$ rad/s (equivalent to $f_0 = 18.96$ Hz). The external force $\boldsymbol{F}$ is a single point vertical load.

Damage in the material is modelled through a damage or phase-field variable denoted by $d$, following the hybrid approach in Ambati et al. (2015). The damage field takes value 0 at intact points of the material and value 1 at fully broken parts of it, with a smooth transition between the two states. The stress-strain constitutive relation in (1) is degraded by a function $g(d)$ to account for the loss of stiffness in damaged regions. We assume quadratic degradation

$$g(d) := \begin{cases} (1 - d)^2 & \text{where } \Psi_0^+ \geq \Psi_0^-, \\ 1 & \text{otherwise.} \end{cases} \tag{2}$$

Here $\Psi_0^+$ denotes the tensile component of the elastic energy density, in opposition to the compressive component $\Psi_0^-$. The condition in (2) ensures that no interpenetration of faces occurs under compression, restoring the original stiffness of the material when compression dominates over tension. The condition is evaluated using the displacement field at the previous time step. This is an alternative to incorporating the splitting in (1) that allows to maintain a linear equilibrium equation. Here, we use the tension-compression splitting proposed in Amor, Marigo and Maurini (2009),

$$\Psi_0^+(\boldsymbol{\varepsilon}) = \frac{1}{2}K\langle \text{tr}(\boldsymbol{\varepsilon})\rangle_+^2 + \mu\left(\boldsymbol{\varepsilon}^{dev} : \boldsymbol{\varepsilon}^{dev}\right), \quad \text{and} \quad \Psi_0^-(\boldsymbol{\varepsilon}) = \frac{1}{2}K\langle \text{tr}(\boldsymbol{\varepsilon})\rangle_-^2,$$

with $\langle a \rangle_\pm = \frac{1}{2}(a \pm |a|)$, $K$ the bulk modulus and $\boldsymbol{\varepsilon}^{dev}$ the deviatoric part of $\boldsymbol{\varepsilon}$.

For the undamaged beam, $d = 0$ in the entire domain. In the faulty configurations, the notch is introduced into the domain through a damage field $d$, which is computed following the approach of Borden, Verhoosel, Scott, Hughes and Landis (2012) with phase-field regularization length $\ell = 0.02$ m. We account for different cases with 36% loss of stiffness across the notch (corresponding to $d = 0.2$) and a complete loss of stiffness ($d = 1$). Note that, differently to standard phase-field simulations, no evolution of the damage is considered, *i.e.*, $d$ is fixed for all time. Also, no artificial stiffness is introduced in (2).

The model in (1) is solved using a linear Finite Element discretization in space and the generalized-$\alpha$ method to integrate in time Chung and Hulbert (1993). The spectral radius for the high-frequency dissipation is set to $\rho_\infty = 0.5$. The method is unconditionally stable, has second-order accuracy in time for displacements and velocities, and optimally balances the dissipation of high and low frequencies.

For the undamaged beam, the computational mesh is uniform and has 192 triangular elements (125 nodes). For the faulty cases, the mesh is refined near the notch to properly capture the phase-field, with an element size $h = \ell/2$. The refinement results in a mesh with 1960 elements in the cases with notch length $L = 0.2$ m (1020 nodes), and with 5078 elements when $L = 0.7$ m (2581 nodes). All computations are performed using the open-source library FeniCS Langtangen and Logg (2017).