

Supporting Material

ICON: an adaptation of infinite HMMs for time traces with drift

I. Sgouralis and S. Pressé

In this Supplement we provide: (i) a description of the source code and GUI implementing ICON as describe in the main text, Sect. S1; (ii) additional time series analyses demonstrating the effects of the drift function $f(t)$, Sect. S2; and (iii) a complete summary of the involved models throughout the main text, Sect. S3.

S1 ICON implementation

The implementation, graphical user interface (GUI), and source code of ICON is similar to the implementation of the plain iHMM described in the companion perspectives article. Here we highlight only the main differences. For further details and an in depth presentation of the offered capabilities we refer to the documentation found in the perspectives.

The present implementation comes in two versions: *single trace ICON*, and *double trace ICON*. These versions are described in the Methods section of the main text and can be used to analyze single molecule data from one or two traces, respectively.

Additionally to the hyperparameters associated with the emission model described in the perspectives article, both ICON versions require M , the total number of interpolation nodes, to be specified by the user.

In the *single trace ICON* the time series to be analyzed must be provided as an 1D array with one value per time level, while in the *double trace ICON* the time series to be analyzed must be provided as a 2D array with two values per time level.

Single trace ICON exports (i) one drift time series, (ii) one emission trajectory, (iii) one state transition matrix, and (iv) the corresponding individual samples. *Double trace ICON* exports (i) two drift time series, (ii) two emission trajectories, (iii) one state transition matrix, and (iv) the corresponding individual samples.

S2 Additional analyses

Here we demonstrate the effects of the various choices of the interpolation $f(t)$, see Methods section in the main text. Specifically, we consider the following choices: basis functions and number of nodes for the interpolation.

Figure S1 shows the resulting state and drift trajectories for two different choices of bases functions: (i) cubic splines (as in the main text) and (ii) linear polynomials. As can be seen, the resulting estimates are minimally affected by this choices.

Figure ?? shows the resulting state and drift trajectories for three different numbers of nodes: (i) $M = 5$, (ii) $M = 10$ (as in the main text), and (iii) $M = 20$. As can be seen, the resulting estimates are also minimally affected by this choices.

S3 Summary of the models used

Here we present a detailed list of the statistical models presented in this study. For definitions and notation we refer to Section 2 of the main text.

S3.1 Hidden Markov Model (HMM)

$$s_0 = \sigma_0 \tag{1}$$

$$s_n | s_{n-1} \sim \text{Cat}(\tilde{\pi}_{s_{n-1}}), \quad n = 1, 2, \dots, N \tag{2}$$

$$x_n | s_n \sim F_{s_n}, \quad n = 1, 2, \dots, N \tag{3}$$

where $\text{Cat}(\tilde{\pi}_{s_{n-1}})$ is the categorical probability distribution with parameters gathered in the probability vector $\tilde{\pi}_{s_{n-1}} = (\pi_{s_{n-1} \rightarrow \sigma_1}, \pi_{s_{n-1} \rightarrow \sigma_2}, \dots)$.

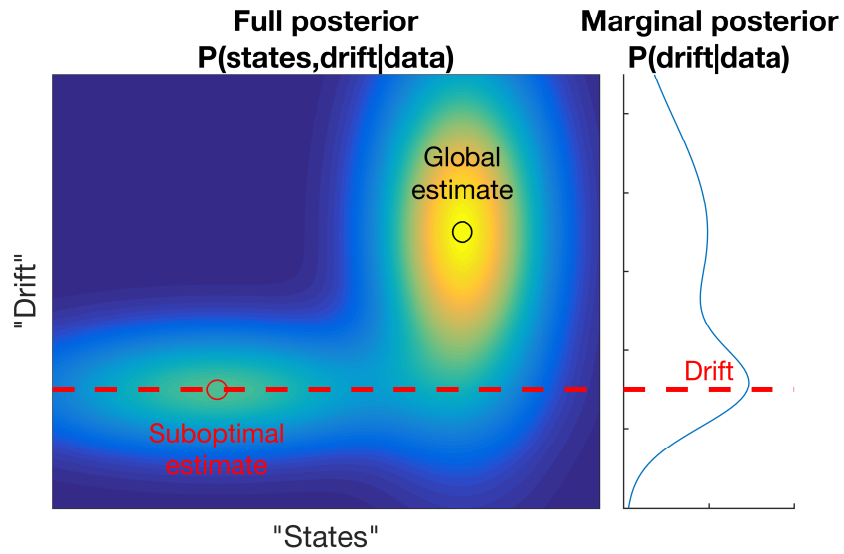


Figure S1: State (left) and drift (right) trajectories estimated using cubic splines and linear polynomials for the interpolation function $f(t)$.

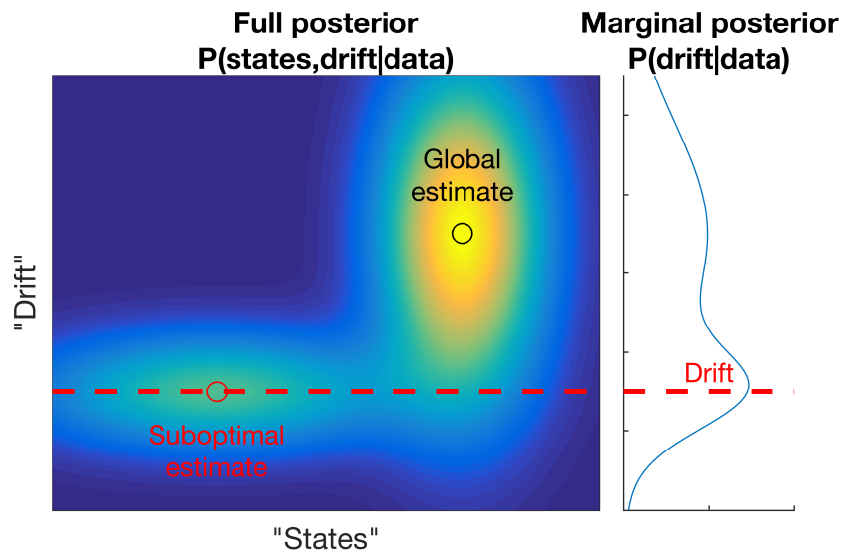


Figure S2: State (left) and drift (right) trajectories estimated using $M = 5, 10$ and 20 in the interpolation function $f(t)$.

S3.2 infinite Hidden Markov Model (iHMM)

$$\tilde{\beta} \sim GEM(\gamma) \quad (4)$$

$$\tilde{\pi}_{\sigma_k} | \tilde{\beta} \sim DP(\alpha, \tilde{\beta}), \quad k = 0, 1, 2, \dots \quad (5)$$

$$\phi_{\sigma_k} \sim H, \quad k = 1, 2, \dots \quad (6)$$

$$s_0 = \sigma_0 \quad (7)$$

$$s_n | s_{n-1}, \tilde{\pi}_{s_{n-1}} \sim Cat(\tilde{\pi}_{s_{n-1}}), \quad n = 1, 2, \dots, N \quad (8)$$

$$x_n | s_n, \phi_{s_n} \sim F_{s_n}, \quad n = 1, 2, \dots, N \quad (9)$$

where $GEM(\gamma)$ is the Griffiths-Engen-McCloskey "stick-breaking" process, $DP(\alpha, \tilde{\beta})$ is the Dirichlet distribution and $\tilde{\beta} = (\beta_{\sigma_1}, \beta_{\sigma_2}, \dots)$ is the base measure.

S3.3 Single trace ICON

$$\tilde{\beta} \sim GEM(\gamma) \quad (10)$$

$$\tilde{\pi}_{\sigma_k} | \tilde{\beta} \sim DP(\alpha, \tilde{\beta}), \quad k = 0, 1, 2, \dots \quad (11)$$

$$\phi_{\sigma_k} \sim H, \quad k = 1, 2, \dots \quad (12)$$

$$s_0 = \sigma_0 \quad (13)$$

$$s_n | s_{n-1}, \tilde{\pi}_{s_{n-1}} \sim Cat(\tilde{\pi}_{s_{n-1}}), \quad n = 1, 2, \dots, N \quad (14)$$

$$x_n | s_n, \phi_{s_n} \sim F_{s_n}, \quad n = 1, 2, \dots, N \quad (15)$$

$$h^* \sim \mathcal{N}(\mu^*, \tau^*) \quad (16)$$

$$w^* \sim \mathcal{G}(a^*, b^*) \quad (17)$$

$$y_m^* | h^*, w^* \sim \mathcal{N}(h^*, w^*), \quad m = 1, 2, \dots, M \quad (18)$$

$$v^* | y_1^*, y_2^*, \dots, y_M^* \sim \delta \quad (19)$$

$$y_n = f(y_1^*, y_2^*, \dots, y_M^*; t_n), \quad n = 1, 2, \dots, N \quad (20)$$

$$z_n = x_n + y_n, \quad n = 1, 2, \dots, N \quad (21)$$

where δ is the Dirac delta, $\mathcal{N}(\mu^*, \tau^*)$ is the normal probability distribution with mean value μ^* and precision τ^* , $\mathcal{G}(a^*, b^*)$ is the gamma probability distribution with shape a^* and scale b^* , and $f(y_1^*, y_2^*, \dots, y_M^*; t)$ is the interpolant of the nodes (t_m^*, y_m^*) evaluated at t .

A graphical model summarizing this model is shown in Fig. S3.

S3.4 Double trace ICON

$$\tilde{\beta} \sim GEM(\gamma) \quad (22)$$

$$\tilde{\pi}_{\sigma_k} | \tilde{\beta} \sim DP(\alpha, \tilde{\beta}), \quad k = 0, 1, 2, \dots \quad (23)$$

$$\phi_{\sigma_k}^1 \sim H^1, \quad k = 1, 2, \dots \quad (24)$$

$$\phi_{\sigma_k}^2 \sim H^2, \quad k = 1, 2, \dots \quad (25)$$

$$s_0 = \sigma_0 \quad (26)$$

$$s_n | s_{n-1}, \tilde{\pi}_{s_{n-1}} \sim Cat(\tilde{\pi}_{s_{n-1}}), \quad n = 1, 2, \dots, N \quad (27)$$

$$x_n^1 | s_n, \phi_{s_n}^1 \sim F_{s_n}^1, \quad n = 1, 2, \dots, N \quad (28)$$

$$x_n^2 | s_n, \phi_{s_n}^2 \sim F_{s_n}^2, \quad n = 1, 2, \dots, N \quad (29)$$

$$h^{*1} \sim \mathcal{N}(\mu^{*1}, \tau^{*1}) \quad (30)$$

$$h^{*2} \sim \mathcal{N}(\mu^{*2}, \tau^{*2}) \quad (31)$$

$$w^{*1} \sim \mathcal{G}(a^{*1}, b^{*1}) \quad (32)$$

$$w^{*2} \sim \mathcal{G}(a^{*2}, b^{*2}) \quad (33)$$

$$y_m^{*1} | h^{*1}, w^{*1} \sim \mathcal{N}(h^{*1}, w^{*1}), \quad m = 1, 2, \dots, M^1 \quad (34)$$

$$y_m^{*2} | h^{*2}, w^{*2} \sim \mathcal{N}(h^{*2}, w^{*2}), \quad m = 1, 2, \dots, M^2 \quad (35)$$

$$v^* | y_1^{*1}, y_2^{*1}, \dots, y_{M^1}^{*1} \sim \delta \quad (36)$$

$$v^{*2} | y_1^{*2}, y_2^{*2}, \dots, y_{M^2}^{*2} \sim \delta \quad (37)$$

$$y_n^1 = f(y_1^{*1}, y_2^{*1}, \dots, y_{M^1}^{*1}; t_n), \quad n = 1, 2, \dots, N \quad (38)$$

$$y_n^2 = f(y_1^{*2}, y_2^{*2}, \dots, y_{M^2}^{*2}; t_n), \quad n = 1, 2, \dots, N \quad (39)$$

$$z_n^1 = x_n^1 + y_n^1, \quad n = 1, 2, \dots, N \quad (40)$$

$$z_n^2 = x_n^2 + y_n^2, \quad n = 1, 2, \dots, N \quad (41)$$

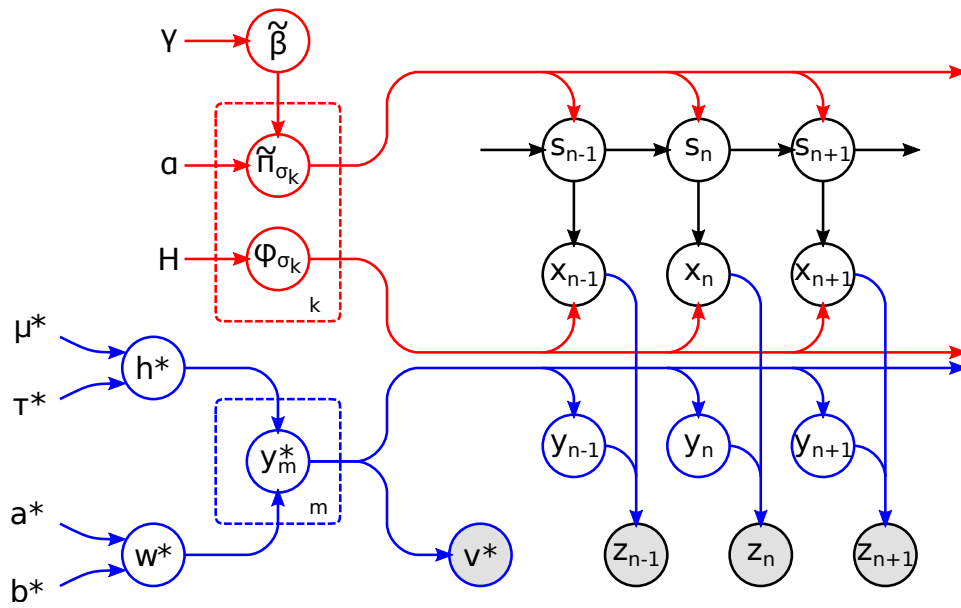


Figure S3: **Graphical representation of single trace ICON.** The hidden Markov model is highlighted with **black**, the prior formulating the infinite hidden Markov model is highlighted with **red**, and the drift representation is highlighted with **blue**.