

Original citation:

Papaspiliopoulos, Omiros, Ruggiero, Matteo and Spanò, Dario. (2016) Conjugacy properties of time-evolving Dirichlet and gamma random measures. *Electronic Journal of Statistics*, 10 (2). pp. 3452-3489.

Permanent WRAP url:

<http://wrap.warwick.ac.uk/83897>

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work of researchers of the University of Warwick available open access under the following conditions.

This article is made available under the Creative Commons Attribution 2.5 Licence (CC BY 2.5) and may be reused according to the conditions of the license. For more details see:

<https://creativecommons.org/licenses/by/2.5/legalcode>

A note on versions:

The version presented in WRAP is the published version, or, version of record, and may be cited as it appears here. For more information, please contact the WRAP Team [at:](mailto:wrap@warwick.ac.uk)

wrap@warwick.ac.uk

Conjugacy properties of time-evolving Dirichlet and gamma random measures

Omiros Papaspiliopoulos*

*ICREA and Department of Economics and Business
Universitat Pompeu Fabra
Ramón Trias Fargas 25-27, 08005, Barcelona, Spain
e-mail: omiros.papaspiliopoulos@upf.edu*

Matteo Ruggiero†

*Collegio Carlo Alberto and ESOMAS Department
University of Torino
C.so Unione Sovietica 218/bis, 10134, Torino, Italy
e-mail: matteo.ruggiero@unito.it*

and

Dario Spanò

*Department of Statistics
University of Warwick
Coventry CV4 7AL, UK
e-mail: D.Spano@warwick.ac.uk*

Abstract: We extend classic characterisations of posterior distributions under Dirichlet process and gamma random measures priors to a dynamic framework. We consider the problem of learning, from indirect observations, two families of time-dependent processes of interest in Bayesian nonparametrics: the first is a dependent Dirichlet process driven by a Fleming–Viot model, and the data are random samples from the process state at discrete times; the second is a collection of dependent gamma random measures driven by a Dawson–Watanabe model, and the data are collected according to a Poisson point process with intensity given by the process state at discrete times. Both driving processes are diffusions taking values in the space of discrete measures whose support varies with time, and are stationary and reversible with respect to Dirichlet and gamma priors respectively. A common methodology is developed to obtain in closed form the time-marginal posteriors given past and present data. These are shown to belong to classes of finite mixtures of Dirichlet processes and gamma random measures for the two models respectively, yielding conjugacy of these classes to the type of data we consider. We provide explicit results on the parameters of the mixture components and on the mixing weights, which are time-varying and drive the mixtures towards the respective priors in absence of further data. Explicit algorithms are provided to recursively compute the parameters of the mixtures. Our results are based on the projective properties of the signals and on certain duality properties of their projections.

*Supported by the MINECO/FEDER via grant MTM2015-67304-P.

†Supported by the European Research Council (ERC) through StG “N-BNP” 306406.

MSC 2010 subject classifications: Primary 62M05, 62M20; secondary 62G05, 60J60, 60G57.

Keywords and phrases: Bayesian nonparametrics, Dawson–Watanabe process, Dirichlet process, duality, Fleming–Viot process, gamma random measure.

Received December 2015.

Contents

1	Introduction	3453
1.1	Motivation and main contributions	3453
1.2	Hidden Markov models	3457
1.3	Illustration for CIR and WF signals	3458
1.4	Preliminary notation	3462
2	Hidden Markov measures	3463
2.1	Fleming–Viot signals	3463
2.1.1	The static model: Dirichlet processes and mixtures thereof	3463
2.1.2	The Fleming–Viot process	3464
2.2	Dawson–Watanabe signals	3466
2.2.1	The static model: Gamma random measures and mixtures thereof	3466
2.2.2	The Dawson–Watanabe process	3467
3	Conjugacy properties of time-evolving Dirichlet and gamma random measures	3468
3.1	Filtering Fleming–Viot signals	3468
3.2	Filtering Dawson–Watanabe signals	3471
4	Theory for computable filtering of FV and DW signals	3474
4.1	Computable filtering and duality	3474
4.2	Computable filtering for Fleming–Viot processes	3477
4.3	Computable filtering for Dawson–Watanabe processes	3480
	Acknowledgements	3486
	References	3486

1. Introduction

1.1. Motivation and main contributions

An active area of research in Bayesian nonparametrics is the construction and the statistical learning of so-called dependent processes. These aim at accommodating weaker forms of dependence than exchangeability among the data, such as partial exchangeability in the sense of de Finetti. The task is then to let the infinite-dimensional parameter, represented by a random measure, depend on a covariate, so that the generated data are exchangeable only conditional on the same covariate value, but not overall exchangeable. This approach was inspired by MacEachern (1999, 2000) and has received considerable attention since.

In the context of this article, the most relevant strand of this literature attempts to build time evolution into standard random measures for semi-parametric time-series analysis, combining the merits of flexible exchangeable modelling afforded by random measures with those of mainstream generalised linear and time series modelling. For the case of Dirichlet processes, the reference model in Bayesian nonparametrics introduced by Ferguson (1973), the time evolution has often been built into the process by exploiting its celebrated stick-breaking representation (Sethuraman, 1994). For example, Dunson (2006) models the dependent process as an autoregression with Dirichlet distributed innovations, Caron et al. (2008) models the noise in a dynamic linear model with a Dirichlet process mixture, Caron et al. (2007) develops a time-varying Dirichlet mixture with reweighing and movement of atoms in the stick-breaking representation, Rodriguez and ter Horst (2008) induces the dependence in time only via the atoms in the stick-breaking representation, by making them into an heteroskedastic random walk. See also Caron and Teh (2012); Caron et al. (2016); Griffin and Steel (2006); Gutierrez et al. (2016); Mena and Ruggiero (2016). The stick-breaking representation of the Dirichlet process has demonstrated its versatility for constructing dependent processes, but makes it hard to derive any analytical information on the posterior structure of the quantities involved. Parallel to these developments, random measures have been combined with hidden Markov time series models, either for allowing the size of the latent space to evolve in time using transitions based on a hierarchy of Dirichlet processes, e.g. Beal et al. (2002); Van Gael et al. (2008); Stepleton et al. (2009); Zhang et al. (2014), or for building flexible emission distributions that link the latent states to data, e.g. Yau et al. (2011); Gassiat and Rousseau (2016).

From a probabilistic perspective, there is a fairly canonical way to build stationary processes with marginal distributions specified as random measures using stochastic differential equations. This more principled approach to building time series with given marginals has been well explored, both probabilistically and statistically, for finite-dimensional marginal distributions, either using processes with discontinuous sample paths, as in Barndorff-Nielsen and Shephard (2001) or Griffin (2011), or using diffusions, as we undertake here. The relevance of measure-valued diffusions in Bayesian nonparametrics has been pioneered in Walker et al. (2007), whose construction naturally allows for separate control of the marginal distributions and the memory of the process.

The statistical models we investigate in this article, introduced in Section 2, can be seen as instances of what we call *hidden Markov measures*, since the models are formulated as hidden Markov models where the latent, unobserved signal is a measure-valued infinite-dimensional Markov process. The signal in the first model is the Fleming–Viot (FV) process, denoted $\{X_t, t \geq 0\}$ on some state space \mathcal{Y} (also called type space in population genetics), which admits the law of a Dirichlet process on \mathcal{Y} as marginal distribution. At times t_n , conditionally on $X_{t_n} = x$, observations are drawn independently from x , i.e.,

$$Y_{t_n,i} \mid x \stackrel{iid}{\sim} x, \quad i = 1, \dots, m_{t_n}, \quad Y_{t_n,i} \in \mathcal{Y}. \quad (1.1)$$

Hence, this statistical model is a dynamic extension of the classic Bayesian non-parametric model for unknown distributions of Ferguson (1973) and Antoniak (1974). The signal in the second model is the Dawson–Watanabe (DW) process, denoted $\{Z_t, t \geq 0\}$ and also defined on \mathcal{Y} , that admits the law of a gamma random measure as marginal distribution. At times t_n , conditionally on $Z_{t_n} = z$, the observations are a Poisson process Y_{t_n} on \mathcal{Y} with random intensity z , i.e., for any collection of disjoint sets $A_1, \dots, A_K \in \mathcal{Y}$ and $K \in \mathbb{N}$,

$$Y_{t_n}(A_i) | z \stackrel{\text{ind}}{\sim} \text{Po}(z(A_i)).$$

Hence, this is a time-evolving Cox process and can be seen as a dynamic extension of the classic Bayesian nonparametric model for Poisson point processes of Lo (1982).

The Dirichlet and the gamma random measures, used as Bayesian nonparametric priors, have conjugacy properties to observation models of the type described above, which have been exploited both for developing theory and for building simulation algorithms for posterior and predictive inference. These properties, reviewed in Sections 2.1.1 and 2.2.1, have propelled the use of these models into mainstream statistics, and have been used directly in simpler models or to build updating equations within Markov chain Monte Carlo and variational Bayes computational algorithms in hierarchical models.

In this article, for the first time, we show that the dynamic versions of these Dirichlet and gamma models also enjoy certain conjugacy properties. First, we formulate such models as hidden Markov models where the latent signal is a measure-valued diffusion and the observations arise at discrete times according to the mechanisms described above. We then obtain that the filtering distributions, that is the laws of the signal at each observation time conditionally on all data up to that time, are finite mixtures of Dirichlet and gamma random measures respectively. We provide a concrete posterior characterisation of these marginal distributions and explicit algorithms for the recursive computation of the parameters of these mixtures. Our results show that these families of finite mixtures are closed with respect to the Bayesian learning in this dynamic framework, and thus provide an extension of the classic posterior characterisations of Antoniak (1974) and Lo (1982) to time-evolving settings.

The techniques we use to establish the new conjugacy results are detailed in Section 4, and build upon three aspects: the characterisations of Dirichlet and gamma random measures through their projections; certain results on measure-valued diffusions related to their time-reversal; and some very recent developments in Papaspiliopoulos and Ruggiero (2014) that relate optimal filtering for finite-dimensional hidden Markov models with the notion of duality for Markov processes, reviewed in Section 4.1. Figure 1 schematises, from a high level perspective, the strategy for obtaining our results. In a nutshell, the essence of our theoretical results is that the operations of projection and propagation of measures commute. More specifically, we first exploit the characterisation of the Dirichlet and gamma random measures via their finite-dimensional distributions, which are Dirichlet and independent gamma distributions respectively.

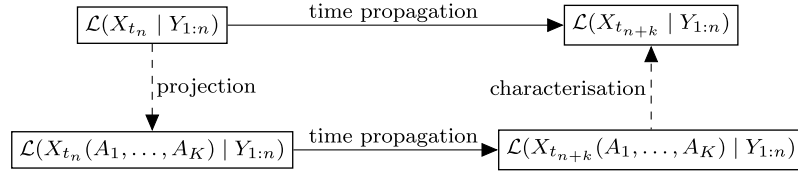


FIG 1. Scheme of the general argument for obtaining the filtering distribution of hidden Markov models with FV and DW signals, proved in Theorems 3.1 and 3.2. In this figure X_t is the latent measure-valued signal. Given data $Y_{1:n}$, the future distribution of the signal $\mathcal{L}(X_{t_{n+k}} | Y_{1:n})$ at time t_{n+k} is determined by taking its finite-dimensional projection $\mathcal{L}(X_{t_n}(A_1, \dots, A_K) | Y_{1:n})$ onto an arbitrary partition (A_1, \dots, A_K) , evaluating the relative propagation $\mathcal{L}(X_{t_{n+k}}(A_1, \dots, A_K) | Y_{1:n})$ at time t_{n+k} , and by exploiting the projective characterisation of the filtering distributions.

Then we exploit the fact that the dynamics of these finite-dimensional distributions induced by the measure-valued signals are the Wright–Fisher (WF) diffusion and a multivariate Cox–Ingersoll–Ross (CIR) diffusion. Then, we extend the results in Papaspiliopoulos and Ruggiero (2014) to show that filtering these finite-dimensional signals on the basis of observations generated as described above results in mixtures of Dirichlet and independent gamma distributions. Finally, we use again the characterisations of Dirichlet and gamma measures via their finite-dimensional distributions to obtain the main results in this paper, that the filtering distributions in the Fleming–Viot model evolves in the family of finite mixtures of Dirichlet processes and those in the Dawson–Watanabe model in the family of finite mixtures of gamma random measures, under the observation models considered. The validity of this argument is formally proved in Theorems 3.1 and 3.2. The resulting recursive procedures for Fleming–Viot and Dawson–Watanabe signals that describe how to compute the parameters of the mixtures at each observation time are given in Propositions 3.1 and 3.2, and the associated pseudo codes are outlined in Algorithms 1 and 2.

The paper is organised as follows. Section 1.2 briefly introduces some basic concepts on hidden Markov models. Section 1.3 provides a simple illustration of the underlying structures implied by previous results on filtering one-dimensional WF and CIR processes. These will be the reference examples throughout the paper and provide relevant intuition on our main results in terms of special cases, since the WF and CIR model are the one-dimensional projections of the infinite-dimensional families we consider here. Section 2 describes the two families of dependent random measures which are the object of this contribution, the Fleming–Viot and the Dawson–Watanabe diffusions, from a non technical viewpoint. Connections of the dynamic models with their marginal or static sub-cases given by Dirichlet and gamma random measures, well known in Bayesian nonparametrics, are emphasised. Section 3 exposes and discusses the main results on the conjugacy properties of the two above families, given observation models as described earlier, together with the implied algorithms for recursive computation. All the technical details related to the strategy for proving the main results and to the duality structures associated to the signals are deferred to Section 4.

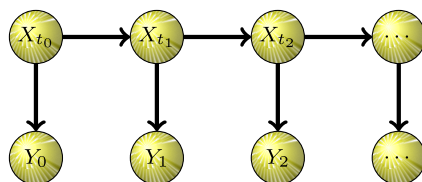


FIG 2. Hidden Markov model represented as a graphical model.

1.2. Hidden Markov models

Since our time-dependent Bayesian nonparametric models are formulated as hidden Markov models, we introduce here some basic related notions. A hidden Markov model (HMM) is a double sequence $\{(X_{t_n}, Y_n), n \geq 0\}$ where X_{t_n} is an unobserved Markov chain, called *latent signal*, and $Y_n := Y_{t_n}$ are conditionally independent observations given the signal. Figure 2 provides a graphical representation of an HMM. We will assume here that the signal is the discrete time sampling of a continuous time Markov process X_t with transition kernel $P_t(x, dx')$. The signal parametrises the law of the observations $\mathcal{L}(Y_n | X_{t_n})$, called *emission distribution*. When this law admits density, this will be denoted by $f_x(y)$.

Filtering optimally an HMM requires the sequential exact evaluation of the so-called *filtering distributions* $\mathcal{L}(X_{t_n} | Y_{0:n})$, i.e., the laws of the signal at different times given past and present observations, where $Y_{0:n} = (Y_1, \dots, Y_n)$. Denote $\nu_n := \mathcal{L}(X_{t_n} | Y_{0:n})$ and let ν be the prior distribution for X_{t_0} . The exact or optimal filter is the solution of the recursion

$$\nu_0 = \phi_{Y_{t_0}}(\nu), \quad \nu_n = \phi_{Y_{t_n}}(\psi_{t_n - t_{n-1}}(\nu_{n-1})), \quad n \in \mathbb{N}. \quad (1.2)$$

This involves the following two operators acting on measures: the *update operator*, which in case a density exists takes the form

$$\phi_y(\nu)(dx) = \frac{f_x(y)\nu(dx)}{p_\nu(y)}, \quad p_\nu(y) = \int_{\mathcal{X}} f_x(y)\nu(dx), \quad (1.3)$$

and the *prediction operator*

$$\psi_t(\nu)(dx') = \int_{\mathcal{X}} \nu(dx) P_t(x, dx'). \quad (1.4)$$

The update operation amounts to an application of Bayes' Theorem to the currently available distribution conditional on the incoming data. The prediction operator propagates forward the current law of the signal of time t according to the transition kernel of the underlying continuous-time latent process. The above recursion (1.2) then alternates updates given the incoming data and predictions of the latent signal as follows:

$$\mathcal{L}(X_{t_0}) \xrightarrow{\text{update}} \mathcal{L}(X_{t_0} | Y_0) \xrightarrow{\text{prediction}} \mathcal{L}(X_{t_1} | Y_0) \xrightarrow{\text{update}} \mathcal{L}(X_{t_1} | Y_0, Y_1) \xrightarrow{\text{prediction}} \dots$$

If X_{t_0} has prior $\nu = \mathcal{L}(X_{t_0})$, then $\nu_0 = \mathcal{L}(X_{t_0}|Y_0)$ is the posterior conditional on the data observed at time t_0 ; ν_1 is the law of the signal at time t_1 obtained by propagating ν_0 of a $t_1 - t_0$ interval and conditioning on the data Y_0, Y_1 observed at time t_0 and t_1 ; and so on.

1.3. Illustration for CIR and WF signals

In order to appreciate the ideas behind the main theoretical results and the Algorithms we develop in this article, we provide some intuition on the corresponding results for one-dimensional hidden Markov models based on Cox–Ingersoll–Ross (CIR) and Wright–Fisher (WF) signals. These are the one-dimensional projections of the DW and FV processes respectively, so informally we could say that a CIR process stands to a DW process as a gamma distribution stands to a gamma random measure, and a one-dimensional WF stands to a FV process as a Beta distribution stands to a Dirichlet process. The results illustrated in this section follow from Papaspiliopoulos and Ruggiero (2014) and are based on the interplay between computable filtering and duality of Markov processes, summarised later in Section 4.1. The developments in this article rely on these results, which are extended to the infinite-dimensional case. Here we highlight the mechanisms underlying the explicit filters with the aid of figures, and postpone the mathematical details to Section 4.

First, let the signal be a one-dimensional Wright–Fisher diffusion on $[0,1]$, with stationary distribution $\pi = \text{Beta}(\alpha, \beta)$ (see Section 2.1.2), which is also taken as the prior ν for the signal at time 0. The signal can be interpreted as the evolving frequency of type-1 individuals in a population of two types whose individuals generate offspring of the same type of the parent, which may be subject to mutation. The observations are assumed to be Bernoulli with success probability given by the signal state. Upon observation of $\mathbf{y}_{t_0} = (y_{t_0,1}, \dots, y_{t_0,m})$, assuming it gives m_1 type-1 and m_2 type-2 individuals with $m = m_1 + m_2$, the prior $\nu = \pi$ is updated as usual via Bayes’ theorem to $\nu_0 = \phi_{\mathbf{y}_{t_0}}(\nu) = \text{Beta}(\alpha + m_1, \beta + m_2)$. Here $\phi_{\mathbf{y}}$ is the update operator (1.3). A forward propagation of these distribution of time t by means of the prediction operator (1.4) yields the finite mixture of Beta distributions

$$\psi_t(\nu_0) = \sum_{(0,0) \leq (i,j) \leq (m_1, m_2)} p_{(m_1, m_2), (i, j)}(t) \text{Beta}(\alpha + i, \beta + j),$$

whose mixing weights depend on t (see Lemma 4.1 below for their precise definition). The propagation of $\text{Beta}(\alpha + m_1, \beta + m_2)$ at time $t_0 + t$ thus yields a mixture of Beta’s with $(m_1 + 1)(m_2 + 1)$ components. The Beta parameters range from $i = m_1, j = m_2$, which represent the full information provided by the collected data, to $i = j = 0$, which represent the null information on the data so that the associated component coincides with the prior. It is useful to identify the indices of the mixture with the nodes of a graph, as in Figure 3-(b), where the red node represent the component with full information, and the yellow nodes the other components, including the prior identified by the origin.

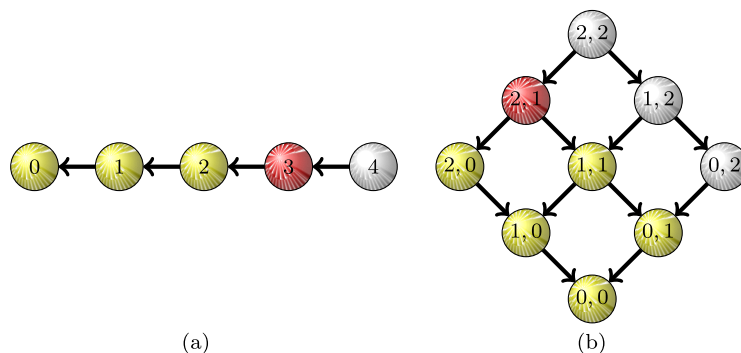


FIG 3. The death process on the lattice modulates the evolution of the mixture weights in the filtering distributions of models with CIR (left) and WF (right) signals. Nodes on the graph identify mixture components in the filtering distribution. The mixture weights are assigned according to the probability that the death process starting from the (red) node which encodes the current full information (here $y = 3$ for the CIR and $(m_1, m_2) = (2, 1)$ for the WF) is in a lower node after time t .

The time-varying mixing weights are the transition probabilities of an associated (dual) 2-dimensional death process, which can be thought of as jumping to lower nodes in the graph of Figure 3-(b) at a specified rate in continuous time. The effect on the mixture of these weights is that as time increases, the probability mass is shifted from components with parameters close to the full information $(\alpha + m_1, \beta + m_2)$, to components which bear less to none information on the data. The mass shift reflects the progressive obsolescence of the data collected at t_0 as evaluated by signal law at time $t_0 + t$ as t increases, and in absence of further data the mixture converges to the prior/stationary distribution π .

Note that it is not obvious that (1.4) yields a finite mixture when P_t is the transition operator of a WF process, since P_t has an infinite series expansion (see Section 2.1.2). This has been proved rather directly in Chaleyat-Maurel and Genon-Catalot (2009) or by combining results on optimal filtering with some duality properties of this model (see Papaspiliopoulos and Ruggiero (2014) or Section 4 here).

Consider now the model where the signal is a one-dimensional CIR diffusion on \mathbb{R}_+ , with gamma stationary distribution (and prior at $t_0 = 0$) given by $\pi = \text{Ga}(\alpha, \beta)$ (see Section 2.2.2). The observations are Poisson with intensity given by the current state of the signal. If the first data are collected at time $t_1 > t_0$, the forward propagation of the signal distribution to time t_1 yields the same distribution by stationarity. Upon observation at time t_1 of $m \geq 1$ Poisson data points with total count y , the prior $\nu = \pi$ is updated via Bayes' theorem to

$$\nu_0 = \text{Ga}(\alpha + y, \beta + m) \quad (1.5)$$

yielding a jump in the measure-valued process; see Figure 4(a). A forward

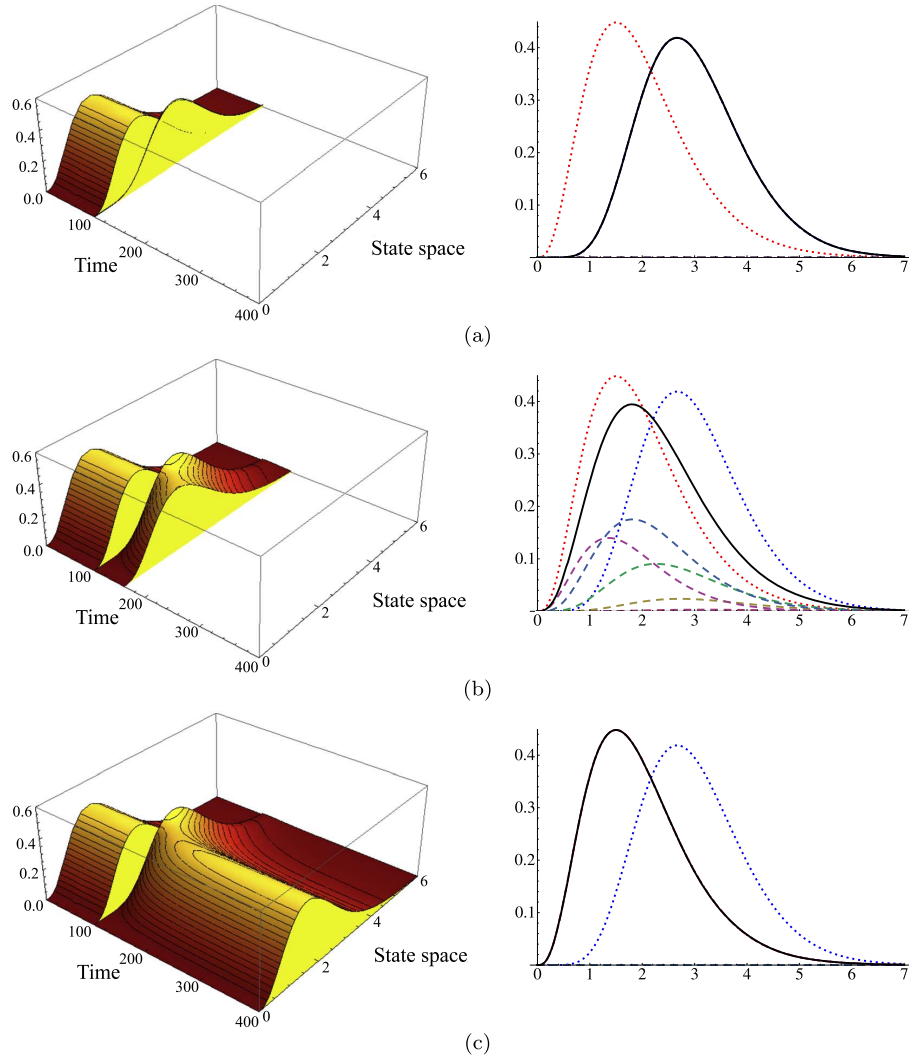


FIG 4. Temporal evolution of the filtering distribution (solid black in right panels and marginal rightmost section of left panels) under the CIR model: (a) until the first data collection the propagation preserves the prior/stationary distribution (red dotted in right panels); at the first data collection, the prior is updated to the posterior (blue dotted in right panels) via Bayes' Theorem, determining a jump in the filtering process (left panel); (b) the forward propagation of the filtering distribution behaves as a finite mixture of Gamma densities (weighted components dashed coloured in right panel); (c) in absence of further data, the time-varying mixing weights shift mass towards the prior component, and the filtering distribution converges to the stationary law.

propagation of ν_0 yields the finite mixture of gamma distributions

$$\psi_t(\nu_0) = \sum_{0 \leq i \leq y} p_{y,i}(t) \text{Ga}(\alpha + i, \beta + S_t), \quad (1.6)$$

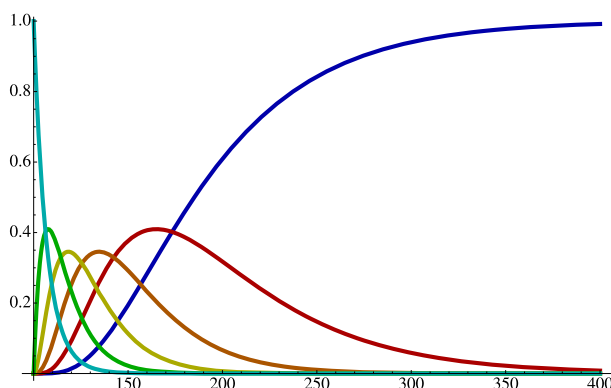


FIG 5. Evolution of the mixture weights which drive the mixture distribution in Fig. 4. At the jump time 100 (the origin here), the mixture component with full posterior information (blue dotted in Fig. 4) has weight equal to 1 (cyan curve), and the other components have zero weight. As the filtering distribution is propagated forward, the weights evolve as transition probabilities of an associated death process. The mixture component equal to the prior distribution (red dotted in Fig. 4), which carries no information on the data, has weight (blue curve) that is 0 at the jump time when the posterior update occurs, and eventually goes back to 1 in absence of further incoming observations, in turn determining the convergence of the mixture to the prior in Fig. 4.

whose mixing weights also depend on t (see Lemma 4.2 below for their precise definition). At time $t_1 + t$, the filtering distribution is a $(y + 1)$ -components mixture with the first gamma parameter ranging from full ($i = y$) to null ($i = 0$) information with respect to the collected data (Figure 4-(b)). The time-dependent mixture weights are the transition probabilities of a certain associated (dual) one-dimensional death process, which can be thought of as jumping to lower nodes in the graph of Figure 3-(a) at a specified rate in continuous time. Similarly to the WF model, the mixing weights shift mass from components whose first parameter is close to the full information, i.e. $(\alpha + y, \beta + S_t)$, to components which bear less to none information $(\alpha, \beta + S_t)$. The time evolution of the mixing weights is depicted in Figure 5, where the cyan and blue lines are the weights of the components with full and no information on the data respectively. As a result of the impact of these weights on the mixture, the latter converges, in absence of further data, to the prior/stationary distribution π as t increases, as shown in Figure 4-(c). Unlike the WF case, in this model there is a second parameter controlled by a deterministic (dual) process S_t on \mathbb{R}_+ which subordinates the transitions of the death process; see Lemma 4.2. Roughly speaking, the death process on the graph controls the obsolescence of the observation counts y , whereas the deterministic process S_t controls that of the sample size m . At the update time t_1 we have $S_0 = m$ as in (1.5), but S_t is a deterministic, continuous and decreasing process, and in absence of further data S_t converges to 0 as $t \rightarrow \infty$, to restore the prior parameter β in the limit of (1.6). See Lemma 4.2 in the Appendix for the formal result for the one-dimensional CIR diffusion.

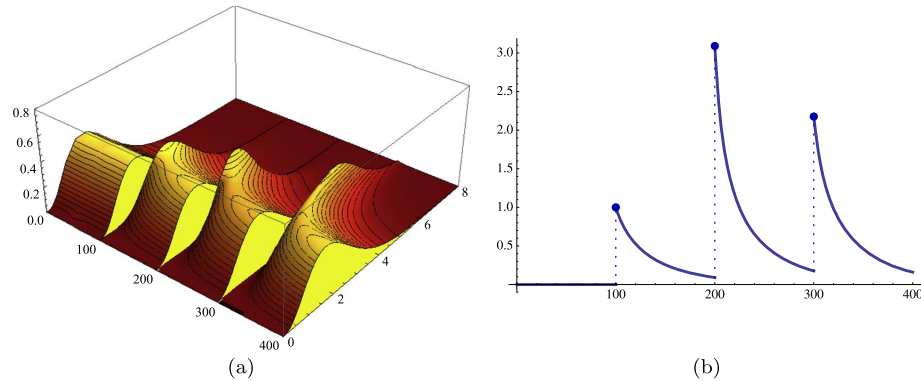


FIG 6. Evolution of the filtering distribution (a) and of the deterministic component of the dual process (b) that modulates the sample size parameter in the mixture components, in the case of multiple data collection at time 100, 200, 300.

When more data samples are collected at different times, the update and propagation operations are alternated, resulting in jump processes for both the filtering distribution and the deterministic dual S_t (Figure 6).

1.4. Preliminary notation

Although most of the notation is better introduced in the appropriate places, we collect here that which is used uniformly over the paper, to avoid recalling these objects several times throughout the text. In all subsequent sections, \mathcal{Y} will denote a locally compact Polish space which represents the observations space, $\mathcal{M}(\mathcal{Y})$ is the associated space of finite Borel measures on \mathcal{Y} and $\mathcal{M}_1(\mathcal{Y})$ its subspace of probability measures. A typical element $\alpha \in \mathcal{M}(\mathcal{Y})$ will be such that

$$\alpha = \theta P_0, \quad \theta > 0, \quad P_0 \in \mathcal{M}_1(\mathcal{Y}), \quad (1.7)$$

where $\theta = \alpha(\mathcal{Y})$ is the total mass of α , and P_0 is sometimes called centering or baseline distribution. We will assume here that P_0 has no atoms. Furthermore, for α as above, Π_α will denote the law on $\mathcal{M}_1(\mathcal{Y})$ of a Dirichlet process, and Γ_α^β that on $\mathcal{M}(\mathcal{Y})$ of a gamma random measure, with $\beta > 0$. These will be recalled formally in Sections 2.1.1 and 2.2.1.

We will denote by X_t the Fleming–Viot process and by Z_t the Dawson–Watanabe process, to be interpreted as $\{X_t, t \geq 0\}$ and $\{Z_t, t \geq 0\}$ when written without argument. Hence X_t and Z_t take values in the space of continuous functions from $[0, \infty)$ to $\mathcal{M}_1(\mathcal{Y})$ and $\mathcal{M}(\mathcal{Y})$ respectively. We will write $X_t(A)$ and $Z_t(A)$ for their respective one dimensional projections onto the Borel set $A \subset \mathcal{Y}$, whereas discrete measures $x(\cdot) \in \mathcal{M}_1(\mathcal{Y})$ and $z(\cdot) \in \mathcal{M}(\mathcal{Y})$ will denote the marginal states of X_t and Z_t . We adopt boldface notation to denote vectors,

with the following conventions:

$$\begin{aligned} \mathbf{x} &= (x_1, \dots, x_K) \in \mathbb{R}_+^K, & \mathbf{m} &= (m_1, \dots, m_K) \in \mathbb{Z}_+^K, \\ \mathbf{x}^{\mathbf{m}} &= x_1^{m_1} \cdots x_K^{m_K}, & |\mathbf{x}| &= \sum_{i=1}^K x_i, \end{aligned}$$

where the dimension $2 \leq K \leq \infty$ will be clear from the context unless specified. Accordingly, the Wright–Fisher model, closely related to projections of the Fleming–Viot process onto partitions, will be denoted \mathbf{X}_t . We denote by $\mathbf{0}$ the vector of zeros and by \mathbf{e}_i the vector whose only non zero entry is a 1 at the i th coordinate. Let also “ $<$ ” define a partial ordering on \mathbb{Z}_+^K , so that $\mathbf{m} < \mathbf{n}$ if $m_j \leq n_j$ for all $j \geq 1$ and $m_j < n_j$ for some $j \geq 1$. Finally, we will use the compact notation $\mathbf{y}_{1:m}$ for vectors of observations y_1, \dots, y_m .

2. Hidden Markov measures

2.1. Fleming–Viot signals

2.1.1. The static model: Dirichlet processes and mixtures thereof

The Dirichlet process on a state space \mathcal{Y} , introduced by Ferguson (1973) (see Ghosal (2010) for a recent review), is a discrete random probability measure $x \in \mathcal{M}_1(\mathcal{Y})$. The process admits the series representation

$$x(\cdot) = \sum_{i=1}^{\infty} W_i \delta_{Y_i}(\cdot), \quad W_i = \frac{Q_i}{\sum_{j \geq 1} Q_j}, \quad Y_i \stackrel{iid}{\sim} P_0, \quad (2.1)$$

where $(Y_i)_{i \geq 1}$ and $(W_i)_{i \geq 1}$ are independent and $(Q_i)_{i \geq 1}$ are the jumps of a gamma process with mean measure $\theta y^{-1} e^{-y} dy$. We will denote by Π_α the law of $x(\cdot)$ in (2.1), with α as in (1.7).

Mixtures of Dirichlet processes were introduced in Antoniak (1974). We say that x is a mixture of Dirichlet processes if

$$x \mid u \sim \Pi_{\alpha_u}, \quad u \sim H,$$

where α_u denotes the measure α conditionally on u , or equivalently

$$x \sim \int_{\mathcal{U}} \Pi_{\alpha_u} dH(u). \quad (2.2)$$

With a slight abuse of terminology we will also refer to the right hand side of the last expression as a mixture of Dirichlet processes.

The Dirichlet process and mixtures thereof have two fundamental properties that are of great interest in statistical learning (Antoniak, 1974):

- *Conjugacy*: let x be as in (2.2). Conditionally on m observations $y_i \mid x \stackrel{iid}{\sim} x$, we have

$$x \mid \mathbf{y}_{1:m} \sim \int_{\mathcal{U}} \Pi_{\alpha_u + \sum_{i=1}^m \delta_{y_i}} dH_{\mathbf{y}_{1:m}}(u),$$

where $H_{\mathbf{y}_{1:m}}$ is the conditional distribution of u given $\mathbf{y}_{1:m}$. Hence a posterior mixture of Dirichlet processes is still a mixture of Dirichlet processes with updated parameters.

- *Projection*: let x be as in (2.2). For any measurable partition A_1, \dots, A_K of \mathcal{Y} , we have

$$(x(A_1), \dots, x(A_K)) \sim \int_{\mathcal{U}} \pi_{\alpha_u} dH(u),$$

where $\alpha_u = (\alpha_u(A_1), \dots, \alpha_u(A_K))$ and π_{α} denotes the Dirichlet distribution with parameter α .

Letting H be concentrated on a single point of \mathcal{U} recovers the respective properties of the Dirichlet process as special case, i.e. $x \sim \Pi_{\alpha}$ and $y_i | x \stackrel{iid}{\sim} x$ imply respectively that $x | \mathbf{y}_{1:m} \sim \Pi_{\alpha + \sum_{i=1}^m \delta_{y_i}}$ and $(x(A_1), \dots, x(A_K)) \sim \pi_{\alpha}$, where $\alpha = (\alpha(A_1), \dots, \alpha(A_K))$.

2.1.2. The Fleming–Viot process

Fleming–Viot (FV) processes are a large family of diffusions taking values in the subspace of $\mathcal{M}_1(\mathcal{Y})$ given by purely atomic probability measures. Hence they describe evolving discrete distributions whose support also varies with time and whose frequencies are each a diffusion on $[0, 1]$. Two states apart in time of a FV process are depicted in Figure 7. See Ethier and Kurtz (1993) and Dawson (1993) for exhaustive reviews. Here we restrict the attention to a subclass known as the (labelled) *infinitely many neutral alleles model* with parent independent mutation, henceforth for simplicity called the FV process, which has the law of a Dirichlet process as stationary measure (Ethier and Kurtz, 1993, Section 9.2).

One of the most intuitive ways to understand a FV process is to consider its transition function, found in Ethier and Griffiths (1993). This is given by

$$P_t(x, dx') = \sum_{m=0}^{\infty} d_m(t) \int_{\mathcal{Y}^m} \Pi_{\alpha + \sum_{i=1}^m \delta_{y_i}}(dx') x^m(dy_1, \dots, dy_m) \quad (2.3)$$

where x^m denotes the m -fold product measure $x \times \dots \times x$ and $\Pi_{\alpha + \sum_{i=1}^m \delta_{y_i}}$ is a posterior Dirichlet process as defined in the previous section. The expression (2.3) has a nice interpretation from the Bayesian learning viewpoint. Given the current state of the process x , with probability $d_m(t)$ an m -sized sample from x is taken, and the arrival state is sampled from the posterior law $\Pi_{\alpha + \sum_{i=1}^m \delta_{y_i}}$. Here $d_m(t)$ is the probability that an N-valued death process which starts at infinity at time 0 is in m at time t , if it jumps from m to $m-1$ at rate $\lambda_m = \frac{1}{2}m(\theta + m - 1)$. See Tavaré (1984) for details. Hence a larger t implies sampling a lower amount of information from x with higher probability, resulting in fewer atoms shared by x and x' . The starting and arrival states thus have correlation which decreases in t as controlled by $d_m(t)$. As $t \rightarrow 0$, infinitely many samples are drawn from x , so x' will coincide with x and the trajectories are continuous in total variation norm (Ethier and Kurtz, 1993). As $t \rightarrow \infty$, the death process which governs

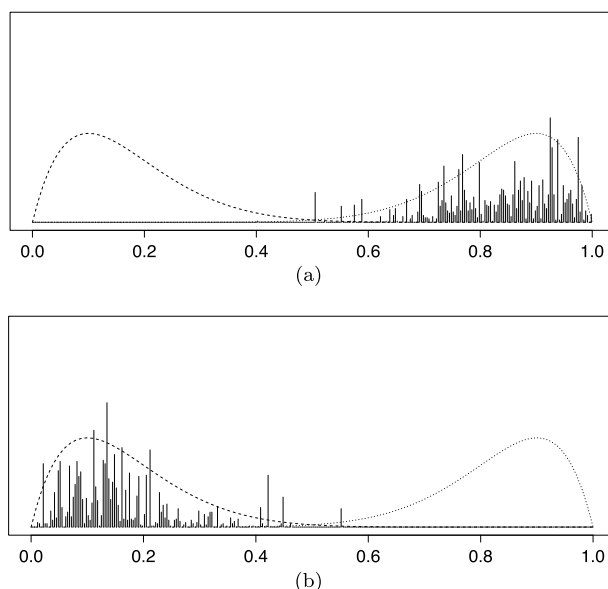


FIG 7. Two states of a FV process on $[0, 1]$ at successive times (solid discrete measures): (a) the initial state has distribution Π_{α_0} with $\alpha_0 = \theta \text{Beta}(4, 2)$ (dotted); (b) after some time, the process reaches the stationary state, which has distribution Π_{α} with $\alpha = \theta \text{Beta}(2, 4)$ (dashed).

the probabilities $d_m(t)$ in (2.3) is eventually absorbed in 0, which implies that $P_t(x, dx') \rightarrow \Pi_{\alpha}$ as $t \rightarrow \infty$, so x' is sampled from the prior Π_{α} . Therefore this FV is stationary with respect to Π_{α} (in fact, it is also reversible). It follows that, using terms familiar to the Bayesian literature, under this parametrisation the FV can be considered as a dependent Dirichlet process with continuous sample paths. Constructions of Fleming–Viot and closely related processes using ideas from Bayesian nonparametrics have been proposed in Walker et al. (2007); Favaro et al. (2009); Ruggiero and Walker (2009a,b). Different classes of diffusive dependent Dirichlet processes or related constructions based on the stick-breaking representation (Sethuraman, 1994) are proposed in Mena and Ruggiero (2016); Mena et al. (2011).

Projecting a FV process X_t onto a measurable partition A_1, \dots, A_K of \mathcal{Y} yields a K -dimensional Wright–Fisher (WF) diffusion \mathbf{X}_t , which is reversible and stationary with respect to the Dirichlet distribution π_{α} , for $\alpha_i = \theta P_0(A_i)$, $i = 1, \dots, K$. See Dawson (2010); Etheridge (2009). This property is the dynamic counterpart of the projective property of Dirichlet processes discussed in Section 2.1.1. Consistently, the transition function of a WF process is obtained as a specialisation of the FV case, yielding

$$P_t(\mathbf{x}, d\mathbf{x}') = \sum_{m=0}^{\infty} d_m(t) \sum_{\mathbf{m} \in \mathbb{Z}_+^K : |\mathbf{m}|=m} \binom{m}{\mathbf{m}} \mathbf{x}^{\mathbf{m}} \pi_{\alpha+\mathbf{m}}(d\mathbf{x}') \quad (2.4)$$

with analogous interpretation to (2.3). See Ethier and Griffiths (1993).

For statistical modelling it is useful to introduce a further parameter σ that controls the speed of the process. This can be done by defining the time change $X_{\tau(t)}$ with $\tau(t) = \sigma t$. In such parameterisation, σ does not affect the stationary distribution of the process, and can be used to model the dependence structure.

2.2. Dawson–Watanabe signals

2.2.1. The static model: Gamma random measures and mixtures thereof

Gamma random measures (Lo, 1982) can be thought of as the counterpart of Dirichlet processes in the context of finite intensity measures. A gamma random measure $z \in \mathcal{M}(\mathcal{Y})$ with shape parameter α as in (1.7) and rate parameter $\beta > 0$, denoted $z \sim \Gamma_{\alpha}^{\beta}$, admits representation

$$z(\cdot) = \beta^{-1} \sum_{i=1}^{\infty} Q_i \delta_{Y_i}(\cdot), \quad Y_i \stackrel{iid}{\sim} P_0, \quad (2.5)$$

with $(Q_i)_{i \geq 1}$ as in (2.1).

Similarly to the definition of mixtures of Dirichlet processes (Section 2.1.1), we say that z is a mixture of gamma random measures if $z \sim \int_{\mathcal{U}} \Gamma_{\alpha_u}^{\beta} dH(u)$, and with a slight abuse of terminology we will also refer to the right hand side of the last expression as a mixture of gamma random measures. Analogous conjugacy and projection properties to those seen for mixtures of Dirichlet processes hold for mixtures of gamma random measures:

- *Conjugacy*: let N be a Poisson point process on \mathcal{Y} with random intensity measure z , i.e., conditionally on z , $N(A_i) \stackrel{ind}{\sim} \text{Po}(z(A_i))$ for any measurable partition A_1, \dots, A_K of \mathcal{Y} , $K \in \mathbb{N}$. Let $m := N(\mathcal{Y})$, and given m , let y_1, \dots, y_m be a realisation of points of N , so that

$$y_i \mid z, m \stackrel{iid}{\sim} z/|z|, \quad m \mid z \sim \text{Po}(|z|) \quad (2.6)$$

where $|z| := z(\mathcal{Y})$ is the total mass of z . Then

$$z \mid \mathbf{y}_{1:m} \sim \int_{\mathcal{U}} \Gamma_{\alpha_u + \sum_{i=1}^m \delta_{y_i}}^{\beta+1} dH_{\mathbf{y}_{1:m}}(u), \quad (2.7)$$

where $H_{\mathbf{y}_{1:m}}$ is the conditional distribution of u given $\mathbf{y}_{1:m}$. Hence mixtures of gamma random measures are conjugate with respect to Poisson point process data.

- *Projection*: for any measurable partition A_1, \dots, A_K of \mathcal{Y} , we have

$$(z(A_1), \dots, z(A_K)) \sim \int_{\mathcal{U}} \prod_{i=1}^K \text{Ga}(\alpha_{u,i}, \beta) dH(u),$$

where $\alpha_{u,i} = \alpha_u(A_i)$, and $\text{Ga}(\alpha, \beta)$ denotes the gamma distribution with shape α and rate β .

Letting H be concentrated on a single point of \mathcal{U} recovers the respective properties of gamma random measures as special case, i.e. $z \sim \Gamma_{\alpha}^{\beta}$ and y_i as in (2.6) imply $z|y_{1:m} \sim \Gamma_{\alpha + \sum_{i=1}^m \delta_{y_i}}^{\beta+1}$, and the vector $(z(A_1), \dots, z(A_K))$ has independent components $z(A_i)$ with gamma distribution $\text{Ga}(\alpha_i, \beta)$, $\alpha_i = \alpha(A_i)$.

Finally, it is well known that (2.1) and (2.5) satisfy the relation in distribution

$$x(\cdot) \stackrel{d}{=} \frac{z(\cdot)}{z(\mathcal{Y})} \quad (2.8)$$

where x is independent of $z(\mathcal{Y})$. This extends to the infinite dimensional case the well known relationship between beta and gamma random variables. See for example Daley and Vere-Jones (2008), Example 9.1(e). See also Konno and Shiga (1988) for an extension of (2.8) to the dynamic case concerning FV and DW processes, which requires a random time change.

2.2.2. The Dawson–Watanabe process

Dawson–Watanabe (DW) processes can be considered as dependent models for gamma random measures, and are, roughly speaking, the gamma counterpart of FV processes. More formally, they are branching measure-valued diffusions taking values in the space of finite discrete measures. As in the FV case, they describe evolving discrete measures whose support varies with time and whose masses are each a positive diffusion, but relaxing the constraint of their masses summing to one to that of summing to a finite quantity. See Dawson (1993) and Li (2011) for reviews. Here we are interested in the special case of subcritical branching with immigration, where subcriticality refers to the fact that in the underlying branching population, which can be used to construct the process, the mean number of offspring per individual is less than one. Specifically, we will consider DW processes with transition function

$$P_t(z, dz') = \sum_{m=0}^{\infty} d_m^{|z|, \beta}(t) \int_{\mathcal{Y}^m} \Gamma_{\alpha + \sum_{i=1}^m \delta_{y_i}}^{\beta + S_t^*}(dz')(z/|z|)^m(dy_1, \dots, dy_m). \quad (2.9)$$

where

$$d_m^{|z|, \beta}(t) = \text{Po}\left(m \mid \frac{|z|^{\beta}}{e^{\beta t/2} - 1}\right) \quad \text{and} \quad S_t^* := \frac{\beta}{e^{\beta t/2} - 1}.$$

See Ethier and Griffiths (1993b). The interpretation of (2.9) is similar to that of (2.3): conditional on the current state given by the measure z , m iid samples are drawn from the normalised measure $z/|z|$ and the arrival state z' is sampled from $\Gamma_{\alpha + \sum_{i=1}^m \delta_{y_i}}^{\beta + S_t^*}$. Here the main structural difference with respect to (2.3), apart from the different distributions involved, is that since in general S_t^* is not an integer quantity, the interpretation as sampling the arrival state z' from a posterior gamma law is not formally correct; cf. (2.7). The sample size m is chosen with probability $d_m^{|z|, \beta}(t)$, which is the probability that an \mathbb{N} -valued death process which starts at infinity at time 0 is in m at time t , if it jumps

from m to $m - 1$ at rate $(m\beta/2)(1 - e^{\beta t/2})^{-1}$. See Ethier and Griffiths (1993b) for details. So z and z' will share fewer atoms the farther they are apart in time. The DW process with the above transition is known to be stationary and reversible with respect to the law Γ_α^β of a gamma random measure; cf. (2.5). See Shiga (1990); Ethier and Griffiths (1993b). The Dawson–Watanabe process has been recently considered as a basis to build time-dependent gamma process priors with Markovian evolution in Caron and Teh (2012) and Spanò and Lijoi (2016).

The DW process satisfies a projective property similar to that seen in Section 2.1.2 for the FV process. Let Z_t have transition (2.9). Given a measurable partition A_1, \dots, A_K of \mathcal{Y} , the vector $(Z_t(A_1), \dots, Z_t(A_K))$ has independent components $z_{t,i} = Z_t(A_i)$ each driven by a Cox–Ingersoll–Ross (CIR) diffusion (Cox et al., 1985). These are also subcritical continuous-state branching processes with immigration, reversible and ergodic with respect to a $\text{Ga}(\alpha_i, \beta)$ distribution, with transition function

$$P_t(z_i, dz'_i) = \sum_{m_i=0}^{\infty} \text{Po}\left(m_i \mid \frac{z_i \beta}{e^{\beta t/2} - 1}\right) \text{Ga}\left(dz'_i \mid \alpha_i + m_i, \beta + S_t^*\right). \quad (2.10)$$

As for FV and WF processes, a further parameter σ that controls the speed of the process can be introduced without affecting the stationary distribution. This can be done by defining an appropriate time change that can be used to model the dependence structure.

3. Conjugacy properties of time-evolving Dirichlet and gamma random measures

3.1. Filtering Fleming–Viot signals

Let the latent signal X_t be a FV process with transition function (2.3). We assume that, given the signal state, observations are drawn independently from x , i.e. as in (1.1) with $X_t = x$. Since x is almost surely discrete (Blackwell, 1973), a sample $\mathbf{y}_{1:m} = (y_1, \dots, y_m)$ from x will feature $K_m \leq m$ ties among the observations with positive probability. Denote by $(y_1^*, \dots, y_{K_m}^*)$ the distinct values in $\mathbf{y}_{1:m}$ and by $\mathbf{m} = (m_1, \dots, m_{K_m})$ the associated multiplicities, so that $|\mathbf{m}| = m$. When an additional sample $\mathbf{y}_{m+1:m+n}$ with multiplicities \mathbf{n} becomes available, we adopt the convention that \mathbf{n} adds up to the multiplicities of the types already recorded in $\mathbf{y}_{1:m}$, so that the total multiplicities count is

$$\mathbf{m} + \mathbf{n} = (m_1 + n_1, \dots, m_{K_m} + n_{K_m}, n_{K_m+1}, \dots, n_{K_m+n}). \quad (3.1)$$

The following Lemma states in our notation the special case of the conjugacy for mixtures of Dirichlet processes which is of interest here; see Section 2.1.1. To this end, let

$$\mathcal{M} = \{\mathbf{m} = (m_1, \dots, m_K) \in \mathbb{Z}_+^K, K \in \mathbb{N}\} \quad (3.2)$$

be the space of multiplicities of K types, with partial ordering defined as in Section 1.4. Denote also by $\text{PU}_\alpha(\mathbf{y}_{m+1:m+n} \mid \mathbf{y}_{1:m})$ the joint distribution of $\mathbf{y}_{m+1:m+n}$ given $\mathbf{y}_{1:m}$ when the random measure x is marginalised out, which is determined by the Blackwell–MacQueen Pólya urn predictive scheme (Blackwell and MacQueen, 1973)

$$Y_{m+i+1} \mid \mathbf{y}_{1:m+i} \sim \frac{\theta P_0 + \sum_{j=1}^{m+i} \delta_{y_j}}{\theta + m + i}, \quad i = 0, \dots, n-1.$$

Lemma 3.1. *Let $M \subset \mathcal{M}$, α as in (1.7) and x be the mixture of Dirichlet processes*

$$x \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}},$$

with $\sum_{\mathbf{m} \in M} w_{\mathbf{m}} = 1$. Given an additional n -sized sample $\mathbf{y}_{m+1:m+n}$ from x with multiplicities \mathbf{n} , the update operator (1.3) yields

$$\phi_{\mathbf{y}_{m+1:m+n}} \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} \right) = \sum_{\mathbf{m} \in M} \hat{w}_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_{m+n}} (m_i + n_i) \delta_{y_i^*}}, \quad (3.3)$$

where

$$\hat{w}_{\mathbf{m}} \propto w_{\mathbf{m}} \text{PU}_\alpha(\mathbf{y}_{m+1:m+n} \mid \mathbf{y}_{1:m}). \quad (3.4)$$

Here “ \propto ” denotes proportionality. The updated distribution is thus still a mixture of Dirichlet processes with different multiplicities and possibly new atoms in the parameter measures $\alpha + \sum_{i=1}^{K_{m+n}} (m_i + n_i) \delta_{y_i^*}$.

The following Theorem formalises our main result on FV processes, showing that the family of finite mixtures of Dirichlet processes is conjugate with respect to discretely sampled data as in (1.1) with $X_t = x$. For \mathcal{M} as in (3.2), let

$$\begin{aligned} L(\mathbf{m}) &= \{\mathbf{n} \in \mathcal{M} : \mathbf{0} \leq \mathbf{n} \leq \mathbf{m}\}, \quad \mathbf{m} \in \mathcal{M}, \\ L(M) &= \{\mathbf{n} \in \mathcal{M} : \mathbf{0} \leq \mathbf{n} \leq \mathbf{m}, \mathbf{m} \in M\}, \quad M \subset \mathcal{M}, \end{aligned} \quad (3.5)$$

be the set of nonnegative vectors lower than or equal to \mathbf{m} or to those in M respectively, with “ \leq ” defined as in Section 1.4). For example, in Figure 3, $L(3)$ and $L((1, 2))$ are both given by all yellow and red nodes in each case. Let also

$$p(\mathbf{i}; \mathbf{m}, |\mathbf{i}|) = \binom{|\mathbf{m}|}{|\mathbf{i}|}^{-1} \prod_{j \geq 1} \binom{m_j}{i_j} \quad (3.6)$$

be the multivariate hypergeometric probability function, with parameters $(\mathbf{m}, |\mathbf{i}|)$, evaluated at \mathbf{i} .

Theorem 3.1. *Let ψ_t be the prediction operator (1.4) associated to a FV process with transition operator (2.3). Then the prediction operator yields as t -time-ahead propagation the finite mixture of Dirichlet processes*

$$\psi_t \left(\Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} \right) = \sum_{\mathbf{n} \in L(\mathbf{m})} p_{\mathbf{m}, \mathbf{n}}(t) \Pi_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}, \quad (3.7)$$

with $L(\mathbf{m})$ as in (3.5) and where

$$p_{\mathbf{m}, \mathbf{m}-\mathbf{i}}(t) = \begin{cases} e^{-\lambda_{|\mathbf{m}|}t}, & \mathbf{i} = \mathbf{0} \\ C_{|\mathbf{m}|, |\mathbf{m}|-|\mathbf{i}|}(t)p(\mathbf{i}; \mathbf{m}, |\mathbf{i}|), & \mathbf{0} < \mathbf{i} \leq \mathbf{m}, \end{cases} \quad (3.8)$$

with

$$C_{|\mathbf{m}|, |\mathbf{m}|-|\mathbf{i}|}(t) = \left(\prod_{h=0}^{|\mathbf{i}|-1} \lambda_{|\mathbf{m}|-h} \right) (-1)^{|\mathbf{i}|} \sum_{k=0}^{|\mathbf{i}|} \frac{e^{-\lambda_{|\mathbf{m}|-k}t}}{\prod_{0 \leq h \leq |\mathbf{i}|, h \neq k} (\lambda_{|\mathbf{m}|-k} - \lambda_{|\mathbf{m}|-h})},$$

$\lambda_n = n(\theta + n - 1)/2$ and $p(\mathbf{i}; \mathbf{m}, |\mathbf{i}|)$ as in (3.6).

The transition operator of the FV process thus maps a Dirichlet process at time t_0 into a finite mixture of Dirichlet processes at time $t_0 + t$. The mixing weights are the transition probabilities of a death process on the K_m dimensional lattice, with K_m being as in (3.7) the number of distinct values observed in previous data. The result is obtained by means of the argument described in Figure 1, which is based on the property that the operations of propagating and projecting the signal commute. By projecting the current distribution of the signal onto an arbitrary measurable partition, yielding a mixture of Dirichlet distributions, we can exploit the results for finite dimensional WF signals to yield the associated propagation (Papaspiliopoulos and Ruggiero, 2014). The propagation of the original signal is then obtained by means of the characterisation of mixtures of Dirichlet processes via their projections. See Section 4.2 for a proof. In particular, the result shows that under these assumptions, the prediction operation (1.4) with the transition function (2.3) reduces to a finite sum.

Iterating the update and propagation operations provided by Lemma 3.1 and Theorem 3.1 allows to perform sequential Bayesian inference on a hidden signal of FV type by means of a finite computation. Here the finiteness refers to the fact that the infinite dimensionality due to the transition function of the signal is avoided analytically, without resorting to any stochastic truncation method for (2.3), given, e.g., by Walker (2007); Papaspiliopoulos and Roberts (2008), and the computation can be conducted in closed form.

The following Proposition formalises the recursive algorithm that sequentially evaluates the marginal posterior laws $\mathcal{L}(X_{t_n}|Y_{1:n})$ of a partially observed FV process by alternating the update and propagation operations, and identifies the family of distributions which is closed with respect to these operations. Define the family of finite mixtures of Dirichlet processes

$$\mathcal{F}_{\Pi} = \left\{ \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} : M \subset \mathcal{M}, |M| < \infty, w_{\mathbf{m}} \geq 0, \sum_{\mathbf{m} \in M} w_{\mathbf{m}} = 1 \right\},$$

with \mathcal{M} as in (3.2) and for a fixed α as in (1.7). Define also

$$t(\mathbf{y}, \mathbf{m}) = \mathbf{m} + \mathbf{n}, \quad \mathbf{m} \in \mathbb{Z}_+^K$$

so that $t(\mathbf{y}, \mathbf{m})$ is (3.1) if \mathbf{n} are the multiplicities of \mathbf{y} , and

$$t(\mathbf{y}, M) = \{\mathbf{n} : \mathbf{n} = t(\mathbf{y}, \mathbf{m}), \mathbf{m} \in M\}, \quad M \subset \mathcal{M}. \quad (3.9)$$

Proposition 3.1. *Let X_t be a FV process with transition function (2.3) and invariant law Π_α defined as in Section 2.1.1, and suppose data are collected as in (1.1) with $X_t = x$. Then \mathcal{F}_Π is closed under the application of the update and prediction operators (1.3) and (1.4). Specifically,*

$$\phi_{\mathbf{y}_{m+1:m+n}} \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} \right) = \sum_{\mathbf{n} \in t(\mathbf{y}_{m+1:m+n}, M)} \hat{w}_{\mathbf{n}} \Pi_{\alpha + \sum_{i=1}^{K_{m+n}} n_i \delta_{y_i^*}}, \quad (3.10)$$

with $t(\mathbf{y}, M)$ as in (3.9),

$$\hat{w}_{\mathbf{n}} \propto w_{\mathbf{m}} \text{PU}_\alpha(\mathbf{y}_{m+1:m+n} \mid \mathbf{y}_{1:m}) \quad \text{for } \mathbf{n} = t(\mathbf{y}, \mathbf{m}), \quad \sum_{\mathbf{n} \in t(\mathbf{y}, M)} \hat{w}_{\mathbf{n}} = 1,$$

and

$$\psi_t \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} \right) = \sum_{\mathbf{n} \in L(M)} p(M, \mathbf{n}, t) \Pi_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}, \quad (3.11)$$

with

$$p(M, \mathbf{n}, t) = \sum_{\mathbf{m} \in M, \mathbf{m} \geq \mathbf{n}} w_{\mathbf{m}} p_{\mathbf{m}, \mathbf{n}}(t) \quad (3.12)$$

and $p_{\mathbf{m}, \mathbf{n}}(t)$ as in (3.8).

Note that the update operation (3.10) preserves the number of components in the mixture, while the prediction operation (3.7) increases its number. The intuition behind this point is analogous to the illustration in Section 1.3, where the prior (node (0, 0)) is updated to the posterior (node (2, 1)) and propagated into a mixture (coloured nodes), with the obvious difference that here the maximum number of distinct values is unbounded and not fixed.

Algorithm 1 describes in pseudo-code the implementation of the filter for FV processes.

3.2. Filtering Dawson–Watanabe signals

Let now the signal Z_t follow a DW process with transition function (2.9), with invariant measure given by the law Γ_α^β of a gamma random measure; see (2.5). We assume that, given the signal state, observations are drawn from a Poisson point process with intensity z , i.e., as in (2.6) with $Z_t = z$. Analogously to the FV case, since z is almost surely discrete, a sample $\mathbf{y}_{1:m} = (y_1, \dots, y_m)$ from (2.6) will feature $K_m \leq m$ ties among the observations with positive probability. To this end, we adopt the same notation as in Section 3.1.

The following Lemma states in our notation the special case of the conjugacy for mixtures of gamma random measures which is of interest here; see Section 2.2.1.

Algorithm 1: Filtering algorithm for FV signals

Data: $\mathbf{y}_{t_j} = (y_{t_j,1}, \dots, y_{t_j,m_{t_j}})$ at times t_j , $j = 0, \dots, J$, as in (1.1)
 Set prior parameters $\alpha = \theta P_0$, $\theta > 0$, $P_0 \in \mathcal{M}_1(\mathcal{Y})$

Initialise
 $\mathbf{y} \leftarrow \emptyset$, $\mathbf{y}^* = \emptyset$, $m \leftarrow 0$, $\mathbf{m} \leftarrow \mathbf{0}$, $M \leftarrow \{\mathbf{0}\}$, $K_m \leftarrow 0$, $w_{\mathbf{0}} \leftarrow 1$

For $j = 0, \dots, J$
 Compute data summaries
 read data \mathbf{y}_{t_j}
 $m \leftarrow m + \text{card}(\mathbf{y}_{t_j})$
 $\mathbf{y}^* \leftarrow$ distinct values in $\mathbf{y}^* \cup \mathbf{y}_{t_j}$
 $K_m \leftarrow \text{card}(\mathbf{y}^*)$
 Update operation
 for $\mathbf{m} \in M$
 $\mathbf{n} \leftarrow t(\mathbf{y}_{t_j}, \mathbf{m})$
 $w_{\mathbf{n}} \leftarrow w_{\mathbf{m}} \text{PU}_{\alpha}(\mathbf{y}_{t_j} \mid \mathbf{y})$
 $M \leftarrow t(\mathbf{y}_{t_j}, M)$
 for $\mathbf{m} \in M$
 $w_{\mathbf{m}} \leftarrow w_{\mathbf{m}} / \sum_{\ell \in M} w_{\ell}$
 $X_{t_j} \mid \mathbf{y}, \mathbf{y}_{t_j} \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}$
 Propagation operation
 for $\mathbf{n} \in L(M)$
 $w_{\mathbf{n}} \leftarrow p(M, \mathbf{n}, t_{j+1} - t_j)$ as in (3.12)
 $M \leftarrow L(M)$
 $X_{t_{j+1}} \mid \mathbf{y}, \mathbf{y}_{t_j} \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}$
 $\mathbf{y} \leftarrow \mathbf{y} \cup \mathbf{y}_{t_j}$

Lemma 3.2. Let \mathcal{M} be as in (3.2), $M \subset \mathcal{M}$, α as in (1.7) and z be the mixture of gamma random measures

$$z \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+1},$$

with $\sum_{\mathbf{m} \in M} w_{\mathbf{m}} = 1$. Given an additional n -sized sample $\mathbf{y}_{m+1:m+n}$ from z as in (2.6) with multiplicities \mathbf{n} , the update operator (1.3) yields

$$\phi_{\mathbf{y}_{m:m+n}} \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+1} \right) = \sum_{\mathbf{m} \in M} \hat{w}_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_{m+n}} (m_i + n_i) \delta_{y_i^*}}^{\beta+2}, \quad (3.13)$$

with $\hat{w}_{\mathbf{m}}$ as in (3.4).

The updated distribution is thus still a mixture of gamma random measures with updated parameters and the same number of components.

The following Theorem formalises our main result on DW processes, showing that the family of finite mixtures of gamma random measures is conjugate with respect to data as in (2.6) with $Z_t = z$.

Theorem 3.2. Let ψ_t be the prediction operator (1.4) associated to a DW process with transition operator (2.9). Let also $L(M)$ be as in (3.5). Then the prediction operator yields as t -time-ahead propagation the finite mixture of gamma random measures

$$\psi_t \left(\Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+s} \right) = \sum_{\mathbf{n} \in L(\mathbf{m})} \tilde{p}_{\mathbf{m},\mathbf{n}}(t) \Gamma_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}^{\beta+S_t}, \quad (3.14)$$

where

$$\tilde{p}_{\mathbf{m},\mathbf{n}}(t) = \text{Bin}(|\mathbf{m}| - |\mathbf{n}|; |\mathbf{m}|, p(t)) p(\mathbf{n}; \mathbf{m}, |\mathbf{n}|), \quad (3.15)$$

and

$$p(t) = S_t/S_0, \quad S_t = \frac{\beta S_0}{(\beta + S_0)e^{\beta t/2} - S_0}, \quad S_0 = s. \quad (3.16)$$

with $p(\mathbf{n}; \mathbf{m}, |\mathbf{n}|)$ as in (3.6) and $\text{Bin}(|\mathbf{m}| - |\mathbf{n}|; |\mathbf{m}|, p(t))$ denoting a Binomial pmf with parameters $(|\mathbf{m}|, p(t))$ evaluated at $|\mathbf{m}| - |\mathbf{n}|$.

The transition operator of the DW process thus maps a gamma random measure into a finite mixture of gamma random measures. The time-varying mixing weights factorise into the binomial transition probabilities of a one-dimensional death process starting at the total size of previous data $|\mathbf{m}|$ and into a hypergeometric pmf. The intuition is that the death process regulates how many levels down the K_m dimensional lattice are taken, and the hypergeometric probability chooses which admissible path down the graph is chosen given the arrival level. In Figure 3 we would have $K_m = 2$ distinct values with multiplicities $\mathbf{m} = (2, 1)$ and total size $|\mathbf{m}| = 3$. Then, e.g., $\tilde{p}_{(2,1),(1,1)}(t)$, is given by the probability $\text{Bin}(1; 3, p(t))$ that the death process jumps down one level from 3 in time t (Figure 3-(a)), times the probability $p((1, 1); (2, 1), 2)$, conditional on going down one level, of reaching $(1, 1)$ from $(2, 1)$ instead of $(2, 0)$, i.e. of removing one item from the pair and not the singleton observation. The Binomial transition of the one-dimensional death process is subordinated to a deterministic process S_t which modulates the sample size continuously in (3.14), starts at the value $S_0 = s$ (cf. the left hand side of (3.14)) and converges to 0 as $t \rightarrow \infty$.

The result is obtained by means of a similar argument to that used for Theorem (3.1), jointly with the relation (2.8) (which here suffices to be applied at the margin of the process). In particular, we exploit the fact that the projection of a DW process onto an arbitrary partition of the space yields a vector of independent CIR processes. See Section 4.3 for a proof. Analogously to the FV case, the result shows that under the present assumptions, the prediction operation (1.4) with the transition function (2.9) reduces to a finite sum.

The following Proposition formalises the recursive algorithm that evaluates the marginal posterior laws $\mathcal{L}(X_{t_n}|Y_{1:n})$ of a partially observed DW process, allowing to perform sequential Bayesian inference on a hidden signal of DW type by means of a finite computation and within the family of finite mixtures of gamma random measures. Define such family as

$$\mathcal{F}_\Gamma = \left\{ \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+s} : \right.$$

$$s > 0, M \subset \mathcal{M}, |M| < \infty, w_{\mathbf{m}} \geq 0, \sum_{\mathbf{m} \in M} w_{\mathbf{m}} = 1 \Big\},$$

with \mathcal{M} as in (3.2).

Proposition 3.2. *Let Z_t be a DW process with transition function (2.9) and invariant law Γ_{α}^{β} defined as in Section 2.2.1, and suppose data are collected as in (2.6) with $Z_t = z$. Then \mathcal{F}_{Γ} is closed under the application of the update and prediction operators (1.3) and (1.4). Specifically,*

$$\phi_{\mathbf{y}_{m+1:m+n}} \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+s} \right) = \sum_{\mathbf{n} \in t(\mathbf{y}_{m+1:m+n}, M)} \hat{w}_{\mathbf{n}} \Gamma_{\alpha + \sum_{i=1}^{K_{m+n}} n_i \delta_{y_i^*}}^{\beta+s+1}, \quad (3.17)$$

with $t(\mathbf{y}, M)$ as in (3.9), $\hat{w}_{\mathbf{n}}$ as in Proposition 3.1, and

$$\psi_t \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+s} \right) = \sum_{\mathbf{n} \in L(M)} p(M, \mathbf{n}, t) \Gamma_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}^{\beta+S_t}. \quad (3.18)$$

with

$$p(M, \mathbf{n}, t) = \sum_{\mathbf{m} \in M, \mathbf{m} \geq \mathbf{n}} w_{\mathbf{m}} \tilde{p}_{\mathbf{m}, \mathbf{n}}(t) \quad (3.19)$$

and $\tilde{p}_{\mathbf{m}, \mathbf{n}}(t)$ as in (3.15) and S_t as in (3.16).

Algorithm 2 describes in pseudo-code the implementation of the filter for DW processes.

4. Theory for computable filtering of FV and DW signals

4.1. Computable filtering and duality

A filter is said to be computable if the sequence of filtering distributions (the marginal laws of the signal given past and current data) can be characterised by a set of parameters whose computation is achieved at a cost that grows at most polynomially with the number of observations. See, e.g., Chaleyat-Maurel and Genon-Catalot (2006). Special cases of this framework are finite dimensional filters for which the computational cost is linear in the number of observations, the Kalman filter for linear Gaussian HMMs being the reference model in this setting.

Let \mathcal{X} denote the state space of the HMM. Papaspiliopoulos and Ruggiero (2014) showed that the existence of a computable filter can be established if the following structures are embedded in the model:

Conjugacy: there exists a function $h(x, \mathbf{m}, \theta) \geq 0$, where $x \in \mathcal{X}$, $\mathbf{m} \in \mathbb{Z}_+^K$ for some $K \in \mathbb{N}$, and $\theta \in \mathbb{R}^l$ for some $l \in \mathbb{N}$, and functions $t_1(y, \mathbf{m})$ and $t_2(y, \theta)$ such that $\int h(x, \mathbf{m}, \theta) \pi(dx) = 1$, for all \mathbf{m} and θ , and

$$\phi_y(h(x, \mathbf{m}, \theta) \pi(dx)) = h(x, t_1(y, \mathbf{m}), t_2(y, \theta)) \pi(dx).$$

Algorithm 2: Filtering algorithm for DW signals

Data: $(m_{t_j}, \mathbf{y}_{t_j}) = (m_{t_j}, y_{t_j,1}, \dots, y_{t_j, m_{t_j}})$ at times t_j , $j = 0, \dots, J$, as in (2.6)

Set prior parameters $\alpha = \theta P_0$, $\theta > 0$, $P_0 \in \mathcal{M}_1(\mathcal{Y})$, $\beta > 0$

Initialise

└ $\mathbf{y} \leftarrow \emptyset$, $\mathbf{y}^* = \emptyset$, $m \leftarrow 0$, $\mathbf{m} \leftarrow \mathbf{0}$, $M \leftarrow \{\mathbf{0}\}$, $K_m \leftarrow 0$, $w_0 \leftarrow 1$, $s = 0$

For $j = 0, \dots, J$

┌ **Compute data summaries**

└ read data \mathbf{y}_{t_j}

└ $m \leftarrow m + \text{card}(\mathbf{y}_{t_j})$

└ $\mathbf{y}^* \leftarrow \text{distinct values in } \mathbf{y}^* \cup \mathbf{y}_{t_j}$

└ $K_m \leftarrow \text{card}(\mathbf{y}^*)$

┌ **Update operation**

└ **for** $\mathbf{m} \in M$

└ └ $\mathbf{n} \leftarrow t(\mathbf{y}_{t_j}, \mathbf{m})$

└ └ $w_{\mathbf{n}} \leftarrow w_{\mathbf{m}} \text{PU}_{\alpha}(\mathbf{y}_{t_j} \mid \mathbf{y})$

└ $M \leftarrow t(\mathbf{y}_{t_j}, M)$

└ **for** $\mathbf{m} \in M$

└ └ $w_{\mathbf{m}} \leftarrow w_{\mathbf{m}} / \sum_{\ell \in M} w_{\ell}$

└ $X_{t_j} \mid \mathbf{y}, \mathbf{y}_{t_j} \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{\mathbf{y}_i^*}}^{\beta + s}$

┌ **Propagation operation**

└ **for** $\mathbf{n} \in L(M)$

└ └ $w_{\mathbf{n}} \leftarrow p(M, \mathbf{n}, t_{j+1} - t_j)$ as in (3.19)

└ $M \leftarrow L(M)$

└ $s' \leftarrow S_{t_{j+1}-t_j}$ as in (3.16), $S_0 = s$

└ $X_{t_{j+1}} \mid \mathbf{y}, \mathbf{y}_{t_j} \sim \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{\mathbf{y}_i^*}}^{\beta + s'}$

└ $s \leftarrow s'$

└ $\mathbf{y} \leftarrow \mathbf{y} \cup \mathbf{y}_{t_j}$

Here $h(x, \mathbf{m}, \theta)\pi(dx)$ identifies a parametric family of distributions which is closed under Bayesian updating with respect to the observation model. Two types of parameters are considered, a multi-index \mathbf{m} and a vector of real-valued parameters θ . The update operator ϕ_y maps the distribution $h(x, \mathbf{m}, \theta)\pi(dx)$, conditional on the new observation y , into a density of the same family with updated parameters $t_1(y, \mathbf{m})$ and $t_2(y, \theta)$. Typically $\pi(dx)$ is the prior and $h(x, \mathbf{m}, \theta)$ is the Radon–Nikodym derivative of the posterior with respect to the prior, when the model is dominated. See, e.g., (4.6) below for an example of such h when π is the Dirichlet distribution.

Duality: there exists a two-component Markov process (M_t, Θ_t) with state-space $\mathbb{Z}_+^K \times \mathbb{R}^l$ and infinitesimal generator

$$(Ag)(\mathbf{m}, \theta) = \lambda(|\mathbf{m}|)\rho(\theta) \sum_{i=1}^K m_i [g(\mathbf{m} - \mathbf{e}_i, \theta) - g(\mathbf{m}, \theta)] + \sum_{i=1}^l r_i(\theta) \frac{\partial g(\mathbf{m}, \theta)}{\partial \theta}$$

acting on bounded functions, such that (M_t, Θ_t) is *dual* to X_t with respect

to the function h , i.e., it satisfies

$$\mathbb{E}^x[h(X_t, \mathbf{m}, \theta)] = \mathbb{E}^{(\mathbf{m}, \theta)}[h(x, M_t, \Theta_t)], \quad (4.1)$$

for all $x \in \mathcal{X}$, $\mathbf{m} \in \mathbb{Z}_+^K$, $\theta \in \mathbb{R}^l$, $t \geq 0$. Here M_t is a death process on \mathbb{Z}_+^K , i.e. a non-increasing pure-jump continuous time Markov process, which jumps from \mathbf{m} to $\mathbf{m} - \mathbf{e}_i$ at rate $\lambda(|\mathbf{m}|)m_i\rho(\theta)$ and is eventually absorbed at the origin; Θ_t is a deterministic process assumed to evolve autonomously according to a system of ordinary differential equations $r(\Theta_t) = d\Theta_t/dt$ for some initial condition $\Theta_0 = \theta_0$ and a suitable function $r: \mathbb{R}^l \rightarrow \mathbb{R}^l$, whose i th coordinate is denoted by r_i in the generator A above and modulates the death rates of M_t through $\rho(\theta)$. The expectations on the left and right hand sides are taken with respect to the law of X_t and (M_t, Θ_t) respectively, conditional on the respective starting points.

The duality condition (4.1) hides a specific distributional relationship between the signal process X_t , which can be thought of as the forward process, and the dual process (M_t, Θ_t) , which can be thought of as unveiling some features of the time reversal structure of X_t . Informally, the death process can be considered as the time reversal of collecting data points if they come at random times, and the deterministic process, in the CIR example (see Section 1.3), can be considered as a continuous reversal of the sample size process, which instead increases by steps. For example, in the well known duality relation between the WF diffusion and the block counting process of Kingman's coalescent, the latter describes the number of surviving non mutant lines of descent in the tree backwards in time which tracks the ancestors of a sample of individuals in the current population. See Griffiths and Spanò (2010). See also Jansen and Kurt (2014) for a review of duality structures for Markov processes.

Note that a local sufficient condition for (4.1), usually easier to check, is

$$(\mathcal{A}h(\cdot, \mathbf{m}, \theta))(x) = (Ah(x, \cdot, \cdot))(\mathbf{m}, \theta), \quad (4.2)$$

for all $\forall x \in \mathcal{X}$, $\mathbf{m} \in \mathbb{Z}_+^K$, $\theta \in \mathbb{R}^l$, where A is as above and \mathcal{A} denotes the infinitesimal generator of the signal X_t .

Under the above conditions, Proposition 2.3 of Papaspiliopoulos and Ruggiero (2014) shows that given the family of distributions

$$\mathcal{F} = \left\{ h(x, \mathbf{m}, \theta)\pi(dx), \mathbf{m} \in \mathbb{Z}_+^K, \theta \in \mathbb{R}^l \right\},$$

if $\nu \in \mathcal{F}$, then the filtering distribution ν_n which satisfies (1.2) is a finite mixture of distributions in \mathcal{F} with parameters that can be computed recursively. This in turn implies that the family of finite mixtures of elements of \mathcal{F} is closed under the iteration of update and prediction operations.

The interpretation is along the lines of the illustration of Section 1.3. Here π , the stationary measure of the forward process, plays the role of the prior distribution and is represented by the origin of \mathbb{Z}_+^K (see Figure 3), which encodes the lack of information on the data generating distribution. Given a sample from

the conjugate observation model, a single component posterior distribution is identified by a node different from the origin in \mathbb{Z}_+^K . The propagation operator then gives positive mass at all nodes which lie beneath the current nodes with positive mass. By iteration of these operations, the filtering distribution evolves within the family of finite mixtures of elements of \mathcal{F} .

4.2. Computable filtering for Fleming–Viot processes

In the present and the following Section we adopt the same notation used in Section 3. We start by formally stating the precise form for the transition probabilities of the death processes involved in the FV filtering. Here the key point to observe is that since the number of distinct types observed in the discrete samples from a FV process is $K_m \leq m$, we only need to consider a generic death processes on $\mathbb{Z}_+^{K_m}$ and not on \mathbb{Z}_+^∞ . For FV processes, the deterministic component Θ_t is constant: here we set $\Theta_t = 1$ for every t and we omit θ from the arguments of the duality function h .

The following Lemma will provide the building block for the proof of Theorem 3.1. In particular, it shows that the transition probabilities of the dual death process are of the form required as coefficients in the expansion (3.8).

Lemma 4.1. *Let $M_t \subset \mathbb{Z}_+^\infty$ be a death process that starts from $M_0 = \mathbf{m}_0 \in \mathcal{M}$, \mathcal{M} as in (3.2), and jumps from \mathbf{m} to $\mathbf{m} - \mathbf{e}_i$ at rate $m_i(\theta + \mathbf{m} - 1)/2$, with generator*

$$\frac{\theta + |\mathbf{m}| - 1}{2} \sum_{i \geq 1} m_i h(\mathbf{x}, \mathbf{m} - \mathbf{e}_i) - \frac{|\mathbf{m}|(\theta + |\mathbf{m}| - 1)}{2} h(\mathbf{x}, \mathbf{m}).$$

Then the transition probabilities for M_t are

$$p_{\mathbf{m}, \mathbf{m} - \mathbf{i}}(t) = \begin{cases} e^{-\lambda_{|\mathbf{m}|} t}, & \mathbf{i} = \mathbf{0}, \\ C_{|\mathbf{m}|, |\mathbf{m}| - |\mathbf{i}|}(t) p(\mathbf{i}; \mathbf{m}, |\mathbf{i}|), & \mathbf{0} < \mathbf{i} \leq \mathbf{m}, \end{cases} \quad (4.3)$$

where

$$C_{|\mathbf{m}|, |\mathbf{m}| - |\mathbf{i}|}(t) = \left(\prod_{h=0}^{|\mathbf{i}|-1} \lambda_{|\mathbf{m}| - h} \right) (-1)^{|\mathbf{i}|} \sum_{k=0}^{|\mathbf{i}|} \frac{e^{-\lambda_{|\mathbf{m}| - k} t}}{\prod_{0 \leq h \leq |\mathbf{i}|, h \neq k} (\lambda_{|\mathbf{m}| - k} - \lambda_{|\mathbf{m}| - h})},$$

$\lambda_n = n(\theta + n - 1)/2$ and $p(\mathbf{i}; \mathbf{m}, |\mathbf{i}|)$ as in (3.6), and 0 otherwise.

Proof. Since $|\mathbf{m}_0| < \infty$, for any such \mathbf{m}_0 the proof is analogous to that of Proposition 2.1 in Papaspiliopoulos and Ruggiero (2014). \square

The following Proof of the conjugacy for mixtures of Dirichlet processes is due to Antoniak (1974) and outlined here for the ease of the reader.

Proof of Lemma 3.1. The distribution x is a mixture of Dirichlet processes with mixing measure $H(\cdot) = \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \delta_{\mathbf{m}}(\cdot)$ on M and transition measure

$$\alpha_{\mathbf{m}}(\cdot) = \alpha(\cdot) + \sum_{j=1}^{K_m} m_j \delta_{y_j^*}(\cdot) = \alpha(\cdot) + \sum_{i=1}^m \delta_{y_i}(\cdot),$$

where $\mathbf{y}_{1:m}$ is the full sample. See Section 2.1.1. Lemma 1 and Corollary 3.2' in Antoniak (1974) now imply that

$$x \mid \mathbf{m}, \mathbf{y}_{m+1:m+n} \sim \Pi_{\alpha_{\mathbf{m}}(\cdot) + \sum_{i=m+1}^n \delta_{y_i}(\cdot)} = \Pi_{\alpha(\cdot) + \sum_{i=1}^n \delta_{y_i}}$$

and $H(\mathbf{m} \mid \mathbf{y}_{m+1:m+n}) \propto w_{\mathbf{m}} \text{PU}_{\alpha}(\mathbf{y}_{m+1:m+n} \mid \mathbf{y}_{1:m})$. \square

As preparatory for the main result on FV processes, we derive here in detail the propagation step for WF processes, which is due to Papaspiliopoulos and Ruggiero (2014). Let

$$\mathcal{A}_K f(\mathbf{x}) = \frac{1}{2} \sum_{i,j=1}^K x_i (\delta_{ij} - x_j) \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} + \frac{1}{2} \sum_{i=1}^K (\alpha_i - \theta x_i) \frac{\partial f(\mathbf{x})}{\partial x_i} \quad (4.4)$$

be the infinitesimal generator of a K -dimensional WF diffusion, with $\alpha_i > 0$ and $\sum_i \alpha_i = \theta$. Here δ_{ij} denotes Kronecker delta and \mathcal{A}_K acts on $C^2(\Delta_K)$ functions, with

$$\Delta_K = \left\{ x \in [0, 1]^K : \sum_{i=1}^K x_i = 1 \right\}.$$

Proposition 4.1. *Let \mathbf{X}_t be a WF diffusion with generator (4.4) and Dirichlet invariant measure on (4.5) denoted π_{α} . Then, for any $\mathbf{m} \in \mathbb{Z}_+^K$ such that $|\mathbf{m}| < \infty$,*

$$\psi_t(\pi_{\alpha+\mathbf{m}}) = \sum_{\mathbf{0} \leq \mathbf{i} \leq \mathbf{m}} p_{\mathbf{m}, \mathbf{m}-\mathbf{i}}(t) \pi_{\alpha+\mathbf{m}-\mathbf{i}}, \quad (4.5)$$

with $p_{\mathbf{m}, \mathbf{m}-\mathbf{i}}(t)$ as in (4.1).

Proof. Define

$$h(\mathbf{x}, \mathbf{m}) = \frac{\Gamma(\theta + |\mathbf{m}|)}{\Gamma(\theta)} \prod_{i=1}^K \frac{\Gamma(\alpha_i)}{\Gamma(\alpha_i + m_i)} \mathbf{x}^{\mathbf{m}}, \quad (4.6)$$

which is in the domain of \mathcal{A}_K . A direct computation shows that

$$\begin{aligned} \mathcal{A}_K h(\mathbf{x}, \mathbf{m}) &= \sum_{i=1}^K \left(\frac{\alpha_i m_i}{2} + \binom{m_i}{2} \right) \frac{\Gamma(\theta + |\mathbf{m}|)}{\Gamma(\theta)} \prod_{j=1}^K \frac{\Gamma(\alpha_j)}{\Gamma(\alpha_j + m_j)} \mathbf{x}^{\mathbf{m} - \mathbf{e}_i} \\ &\quad - \sum_{i=1}^K \left(\frac{\theta m_i}{2} + \binom{m_i}{2} + \frac{1}{2} m_i \sum_{j \neq i} m_j \right) \frac{\Gamma(\theta + |\mathbf{m}|)}{\Gamma(\theta)} \prod_{j=1}^K \frac{\Gamma(\alpha_j)}{\Gamma(\alpha_j + m_j)} \mathbf{x}^{\mathbf{m}} \\ &= \frac{\theta + |\mathbf{m}| - 1}{2} \sum_{i=1}^K m_i h(\mathbf{x}, \mathbf{m} - \mathbf{e}_i) - \frac{|\mathbf{m}|(\theta + |\mathbf{m}| - 1)}{2} h(\mathbf{x}, \mathbf{m}). \end{aligned}$$

Hence, by (4.2), the death process M_t on \mathbb{Z}_+^K , which jumps from \mathbf{m} to $\mathbf{m} - \mathbf{e}_i$ at rate $m_i(\theta + |\mathbf{m}| - 1)/2$, is dual to the WF diffusion with generator \mathcal{A}_K with

respect to (4.6). From the definition (1.4) of the prediction operator now we have

$$\begin{aligned}
 \psi_t(\pi_{\alpha+\mathbf{m}})(d\mathbf{x}') &= \int_{\mathcal{X}} h(\mathbf{x}, \mathbf{m}) \pi_{\alpha}(d\mathbf{x}) P_t(\mathbf{x}, d\mathbf{x}') \\
 &= \int_{\mathcal{X}} h(\mathbf{x}, \mathbf{m}) \pi_{\alpha}(d\mathbf{x}') P_t(\mathbf{x}', d\mathbf{x}) \\
 &= \pi_{\alpha}(d\mathbf{x}') \mathbb{E}^{\mathbf{x}'}[h(\mathbf{X}_t, \mathbf{m})] \\
 &= \pi_{\alpha}(d\mathbf{x}') \mathbb{E}^{\mathbf{m}}[h(\mathbf{x}', M_t)] \\
 &= \pi_{\alpha}(d\mathbf{x}') \sum_{\mathbf{0} \leq \mathbf{i} \leq \mathbf{m}} p_{\mathbf{m}, \mathbf{m}-\mathbf{i}}(t) h(\mathbf{x}', \mathbf{m} - \mathbf{i}) \\
 &= \sum_{\mathbf{0} \leq \mathbf{i} \leq \mathbf{m}} p_{\mathbf{m}, \mathbf{m}-\mathbf{i}}(t) \pi_{\alpha+\mathbf{m}-\mathbf{i}}(d\mathbf{x}')
 \end{aligned}$$

where the second equality holds in virtue of the reversibility of \mathbf{X}_t with respect to π_{α} , the fourth by the duality (4.1) established above together with (4.3) and the fifth from Lemma 4.1. \square

The following proves the propagation step for FV processes by making use of the previous result and by exploiting the strategy outlined in Figure 1.

Proof of Theorem 3.1. Fix an arbitrary partition (A_1, \dots, A_K) of \mathcal{Y} with K classes, and denote by $\tilde{\mathbf{m}}$ the multiplicities resulting from binning $\mathbf{y}_{1:m}$ into the corresponding cells. Then

$$\Pi_{\alpha+\sum_{i=1}^{K_m} m_i \delta_{y_i^*}}(A_1, \dots, A_K) \sim \pi_{\alpha+\tilde{\mathbf{m}}}, \quad (4.7)$$

where $\Pi_{\alpha+\sum_{i=1}^{K_m} m_i \delta_{y_i^*}}(A_1, \dots, A_K)$ denotes the law $\Pi_{\alpha+\sum_{i=1}^{K_m} m_i \delta_{y_i^*}}(\cdot)$ evaluated on (A_1, \dots, A_K) . Since the projection onto the same partition of the FV process is a K -dimensional WF process (see Section 2.1.2), from Proposition 4.1 we have

$$\psi_t\left(\Pi_{\alpha+\sum_{i=1}^{K_m} m_i \delta_{y_i^*}}(A_1, \dots, A_K)\right) = \psi_t(\pi_{\alpha+\tilde{\mathbf{m}}}) = \sum_{\mathbf{n} \in L(\tilde{\mathbf{m}})} p_{\tilde{\mathbf{m}}, \mathbf{n}}(t) \pi_{\alpha+\mathbf{n}}.$$

Furthermore, since a Dirichlet process is characterised by its finite-dimensional projections, now it suffices to show that

$$\sum_{\mathbf{n} \in L(\mathbf{m})} p_{\mathbf{m}, \mathbf{n}}(t) \Pi_{\alpha+\sum_{i=1}^{K_m} n_i \delta_{y_i^*}}(A_1, \dots, A_K) = \sum_{\mathbf{n} \in L(\tilde{\mathbf{m}})} p_{\tilde{\mathbf{m}}, \mathbf{n}}(t) \pi_{\alpha+\mathbf{n}}$$

so that the operations of propagation and projection commute. Given (4.7), we only need to show that the mixture weights are consistent with respect to fragmentation and merging of classes, that is

$$\sum_{\mathbf{i} \in L(\mathbf{m}): \tilde{\mathbf{i}}=\mathbf{n}} p_{\mathbf{m}, \mathbf{i}}(t) = p_{\tilde{\mathbf{m}}, \mathbf{n}}(t),$$

where $\tilde{\mathbf{i}}$ denotes the projection of \mathbf{i} onto (A_1, \dots, A_K) . Using (3.8), the previous in turn reduces to

$$\sum_{\mathbf{i} \in L(\mathbf{m}): \tilde{\mathbf{i}} = \mathbf{n}} p(\mathbf{i}; \mathbf{m}, m - i) = p(\mathbf{n}; \tilde{\mathbf{m}}, m - n),$$

which holds by the marginalization properties of the multivariate hypergeometric distribution. Cf. Johnson et al. (1997), equation 39.3. \square

The last needed result to obtain the recursive representation of Proposition 3.1 reduces now to a simple sum rearrangement.

Proof of Proposition 3.1. The update operation (3.10) follows directly from Lemma 3.1. The prediction operation (3.11) for elements of \mathcal{F}_Π follows from Theorem 3.1 together with the linearity of (1.4) and a rearrangement of the sums, so that

$$\begin{aligned} \psi_t & \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}} \right) \\ &= \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \sum_{\mathbf{n} \in L(\mathbf{m})} p_{\mathbf{m}, \mathbf{n}}(t) \Pi_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}} \\ &= \sum_{\mathbf{n} \in L(M)} \left(\sum_{\mathbf{m} \in M, \mathbf{m} \geq \mathbf{n}} w_{\mathbf{m}} p_{\mathbf{m}, \mathbf{n}}(t) \right) \Pi_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}. \end{aligned} \quad \square$$

4.3. Computable filtering for Dawson–Watanabe processes

The following Lemma, used later, recalls the propagation step for one dimensional CIR processes.

Lemma 4.2. *Let $Z_{i,t}$ be a CIR process with generator (4.8) and invariant distribution $\text{Ga}(\alpha_i, \beta)$. Then*

$$\psi_t(\text{Ga}(\alpha_i + m, \beta + s)) = \sum_{j=0}^m \text{Bin}(m - j; m, p(t)) \text{Ga}(\alpha_i + m - j, \beta + S_t),$$

where

$$p(t) = S_t / S_0, \quad S_t = \frac{\beta S_0}{(\beta + S_0)e^{\beta t/2} - S_0}, \quad S_0 = s.$$

Proof. It follows from Section 3.1 in Papaspiliopoulos and Ruggiero (2014) by letting $\alpha = \delta/2$, $\beta = \gamma/\sigma^2$ and $S_t = \Theta_t - \beta$. \square

As preparatory for proving the main result on DW processes, assume the signal $\mathbf{Z}_t = (Z_{1,t}, \dots, Z_{K,t})$ is a vector of independent CIR components $Z_{i,t}$ each with generator

$$\mathcal{B}_i f(z_i) = \frac{1}{2}(\alpha_i - \beta z_i) f'(z_i) + \frac{1}{2} z_i f''(z_i), \quad (4.8)$$

acting on $C^2([0, \infty))$ functions which vanish at infinity. See Kawazu and Watanabe (1971). The next proposition identifies the dual process for \mathbf{Z}_t .

Theorem 4.1. Let $Z_{i,t}$, $i = 1, \dots, K$, be independent CIR processes each with generator (4.8) parametrised by (α_i, β) , respectively. For $\boldsymbol{\alpha} \in \mathbb{R}_+^K$ and $\theta = |\boldsymbol{\alpha}|$, define $h_{\alpha_i}^C : \mathbb{R}_+ \times \mathbb{Z}_+ \times \mathbb{R}_+$ as

$$h_{\alpha_i}^C(z, m, s) = \frac{\Gamma(\alpha_i)}{\Gamma(\alpha_i + m)} \left(\frac{\beta + s}{\beta} \right)^{\alpha_i} (\beta + s)^m z^m e^{-sz}.$$

Let also $h^W : \mathbb{R}_+^K \times \mathbb{Z}_+^K$ be as in (4.6) and define $h : \mathbb{R}_+^K \times \mathbb{Z}_+^K \times \mathbb{R}_+$ as

$$h(\mathbf{z}, \mathbf{m}, s) = h_\theta^C(|\mathbf{z}|, |\mathbf{m}|, s) h^W(\mathbf{x}, \mathbf{m}),$$

where $\mathbf{x} = \mathbf{z}/|\mathbf{z}|$. Then the joint process $\{(Z_{1,t}, \dots, Z_{K,t}), t \geq 0\}$ is dual, in the sense of (4.1), to the process $\{(\mathbf{M}_t, S_t), t \geq 0\} \subset \mathbb{Z}_+^K \times \mathbb{R}_+$ with generator

$$\begin{aligned} Bg(\mathbf{m}, s) &= \frac{1}{2} |\mathbf{m}| (\beta + s) \sum_{i=1}^K \frac{m_i}{|\mathbf{m}|} [g(\mathbf{m} - \mathbf{e}_i, s) - g(\mathbf{m}, s)] \\ &\quad - \frac{1}{2} s (\beta + s) \frac{\partial g(\mathbf{m}, s)}{\partial s} \end{aligned} \quad (4.9)$$

with respect to $h(\mathbf{z}, \mathbf{m}, s)$.

Proof. Throughout the proof, for ease of notation we will write h_i^C instead of $h_{\alpha_i}^C$. Note first that for all $\mathbf{m} \in \mathbb{Z}_+^K$ we have

$$\prod_{i=1}^K h_i^C(z_i, m_i, s) = h_\theta^C(|z|, |\mathbf{m}|, s) h^W(\mathbf{x}, \mathbf{m}), \quad (4.10)$$

where $x_i = z_i/|z|$, which follows from direct computation by multiplying and dividing by the correct ratios of gamma functions and by writing $\prod_{i=1}^K z_i^{m_i} = |z|^m \prod_{i=1}^K x_i^{m_i}$. We show the result for $K = 2$, from which the statement for general K case follows easily. From the independence of the CIR processes, the generator $(Z_{1,t}, Z_{2,t})$ applied to the left hand side of (4.10) is

$$(\mathcal{B}_1 + \mathcal{B}_2) h_1^C h_2^C = h_2^C \mathcal{B}_1 h_1^C + h_1^C \mathcal{B}_2 h_2^C. \quad (4.11)$$

A direct computation shows that

$$\begin{aligned} \mathcal{B}_i h_i^C &= \frac{m_i}{2} (\beta + s) h_i^C(z_i, m_i - 1, s) + \frac{s}{2} (\alpha_i + m_i) h_i^C(z_i, m_i + 1, s) \\ &\quad - \frac{1}{2} (s(\alpha_i + m_i) + m_i(\beta + s)) h_i^C(z_i, m_i, s). \end{aligned}$$

Substituting in the right hand side of (4.11) and collecting terms with the same coefficients gives

$$\begin{aligned} &\frac{\beta + s}{2} \left[m_1 h_1^C(z_1, m_1 - 1, s) h_2^C(z_2, m_2, s) + m_2 h_1^C(z_1, m_1, s) h_2^C(z_2, m_2 - 1, s) \right] \\ &+ \frac{s}{2} \left[(\alpha_1 + m_1) h_1^C(z_1, m_1 + 1, s) h_2^C(z_2, m_2, s) \right. \end{aligned}$$

$$\begin{aligned}
& + (\alpha_2 + m_2)h_1^C(z_1, m_1, s)h_2^C(z_2, m_2 + 1, s) \Big] \\
& - \frac{1}{2}(s(\alpha + m) + m(\beta + s))h_1^C(z_1, m_1, s)h_2^C(z_2, m_2, s)
\end{aligned}$$

with $\alpha = \alpha_1 + \alpha_2$ and $m = m_1 + m_2$. From (4.10) we now have

$$\begin{aligned}
& \frac{\beta + s}{2}h_\theta^C(|z|, m - 1, s) \Big[m_1h^W(\mathbf{x}, \mathbf{m} - \mathbf{e}_1, s) + m_2h^W(\mathbf{x}, \mathbf{m} - \mathbf{e}_2, s) \Big] \\
& + \frac{s}{2}h_\theta^C(|z|, m + 1, s) \Big[(\alpha_1 + m_1)h^W(\mathbf{x}, \mathbf{m} + \mathbf{e}_1, s) \\
& \quad + (\alpha_2 + m_2)h^W(\mathbf{x}, \mathbf{m} + \mathbf{e}_2, s) \Big] \\
& - \frac{1}{2}(s(\alpha + m) + m(\beta + s))h_\theta^C(|z|, m, s)h^W(\mathbf{x}, \mathbf{m}, s).
\end{aligned}$$

Then

$$\begin{aligned}
& (\mathcal{B}_1 + \mathcal{B}_2)h_1^Ch_2^C \\
& = \frac{\beta + s}{2} \Big[m_1h(\mathbf{z}, \mathbf{m} - \mathbf{e}_1, s) + m_2h(\mathbf{z}, \mathbf{m} - \mathbf{e}_2, s) \Big] \\
& \quad + \frac{s}{2} \Big[(\alpha_1 + m_1)h(\mathbf{z}, \mathbf{m} + \mathbf{e}_1, s) + (\alpha_2 + m_2)h(\mathbf{z}, \mathbf{m} + \mathbf{e}_2, s) \Big] \\
& \quad - \frac{1}{2}(s(\alpha + m) + m(\beta + s))h(\mathbf{z}, \mathbf{m}, s).
\end{aligned} \tag{4.12}$$

Noting now that

$$\begin{aligned}
\frac{\partial}{\partial s}h(\mathbf{z}, \mathbf{m}, s) & = \frac{\alpha + m}{\beta + s}h(\mathbf{z}, \mathbf{m}, s) - \frac{\alpha_1 + m_1}{\beta + s}h(\mathbf{z}, \mathbf{m} + \mathbf{e}_1, s) \\
& \quad - \frac{\alpha_2 + m_2}{\beta + s}h(\mathbf{z}, \mathbf{m} + \mathbf{e}_2, s),
\end{aligned}$$

an application of (4.9) on $h(\mathbf{z}, \mathbf{m}, s)$ shows that $(Bh(\mathbf{z}, \cdot, \cdot))(\mathbf{m}, s)$ equals the right hand side of (4.12), so that (4.2) holds, giving the result. \square

The previous Theorem extends the gamma-type duality showed for one dimensional CIR processes in Papaspiliopoulos and Ruggiero (2014). Although the components of \mathbf{Z}_t are independent, the result is not entirely trivial. Indeed the one-dimensional CIR process is dual to a two-components process given by a one-dimensional death process and a one-dimensional deterministic dual. The previous result shows that K independent CIR processes have dual not given by a K independent versions of the CIR dual, but by a death process on \mathbb{Z}_+^K modulated by a single deterministic process. Specifically, here the dual component \mathbf{M}_t is a K -dimensional death process on \mathbb{Z}_+^K which, conditionally on S_t , jumps from \mathbf{m} to $\mathbf{m} - \mathbf{e}_i$ at rate $2m_i(\beta + S_t)$, and $S_t \in \mathbb{R}_+$ is a nonnegative deterministic process driven by the logistic type differential equation

$$\frac{dS_t}{dt} = -\frac{1}{2}S_t(\beta + S_t). \tag{4.13}$$

The next Proposition formalises the propagation step for multivariate CIR processes. Denote by $\mathbf{Ga}(\alpha, \beta)$ the product of gamma distributions

$$\text{Ga}(\alpha_1, \beta) \times \cdots \times \text{Ga}(\alpha_K, \beta),$$

with $\alpha = (\alpha_1, \dots, \alpha_K)$.

Proposition 4.2. *Let $\{(Z_{1,t}, \dots, Z_{K,t}), t \geq 0\}$ be as in Theorem 4.1. Then*

$$\begin{aligned} \psi_t(\mathbf{Ga}(\alpha + \mathbf{m}, \beta + s)) &= \\ &= \sum_{i=0}^{|\mathbf{m}|} \text{Bin}(|\mathbf{m}| - i; |\mathbf{m}|, p(t)) \text{Ga}(\theta + |\mathbf{m}| - i, \beta + S_t) \\ &\quad \times \sum_{\mathbf{0} \leq \mathbf{i} \leq \mathbf{m}, |\mathbf{i}|=i} p(\mathbf{i}; \mathbf{m}, i) \pi_{\alpha + \mathbf{m} - \mathbf{i}}, \end{aligned} \quad (4.14)$$

where $\text{Bin}(|\mathbf{m}| - i; |\mathbf{m}|, p(t))$ and $p(\mathbf{i}; \mathbf{m}, i)$ are as in (3.15).

Proof. From independence we have

$$\psi_t(\mathbf{Ga}(\alpha + \mathbf{m}, \beta + s)) = \prod_{i=1}^K \psi_t(\text{Ga}(\alpha_i + m_i, \beta + s)).$$

Using Lemma 4.2 in the Appendix, the previous equals

$$\begin{aligned} &\prod_{i=1}^K \sum_{j=0}^{m_i} \text{Bin}(m_i - j; m_i, p(t)) \text{Ga}(\alpha_i + m_i - j, \beta + S_t) \\ &= \sum_{i_1=0}^{m_1} \text{Bin}(m_1 - i_1; m_1, p(t)) \text{Ga}(\alpha_1 + m_1 - i_1, \beta + S_t) \\ &\quad \times \cdots \times \sum_{i_K=0}^{m_K} \text{Bin}(m_K - i_K; m_K, p(t)) \text{Ga}(\alpha_K + m_K - i_K, \beta + S_t). \end{aligned}$$

Using now the fact that a product of Binomials equals the product of a Binomial and an hypergeometric distribution, we have

$$\sum_{i=0}^{|\mathbf{m}|} \text{Bin}(|\mathbf{m}| - i; |\mathbf{m}|, p(t)) \sum_{\mathbf{0} \leq \mathbf{i} \leq \mathbf{m}, |\mathbf{i}|=i} p(\mathbf{i}; \mathbf{m}, i) \prod_{j=1}^K \text{Ga}(\alpha_j + m_j - i_j, \beta + S_t)$$

which, using (2.8), yields (4.14). Furthermore, (3.16) is obtained by solving (4.13) and by means of the following argument. The one dimensional death process that drives $|\mathbf{M}_t|$ in Theorem 4.1, jumps from $|\mathbf{m}|$ to $|\mathbf{m}| - 1$ at rate $|\mathbf{m}|(\beta + S_t)/2$, see (4.9). The probability that $|\mathbf{M}_t|$ remains in $|\mathbf{m}|$ in $[0, t]$ if it is in $|\mathbf{m}|$ at time 0, here denoted $P(|\mathbf{m}| \mid |\mathbf{m}|, S_t)$, is then

$$P(|\mathbf{m}| \mid |\mathbf{m}|, S_t) = \exp \left\{ -\frac{|\mathbf{m}|}{2} \int_0^t (\beta + S_u) du \right\} = \left(\frac{\beta}{(\beta + s)e^{\beta t/2} - s} \right)^{|\mathbf{m}|}.$$

The probability of a jump from $|\mathbf{m}|$ to $|\mathbf{m}| - 1$ occurring in $[0, t]$ is

$$\begin{aligned}
 & P(|\mathbf{m}| - 1 \mid |\mathbf{m}|, S_t) \\
 &= \int_0^t \exp \left\{ -\frac{|\mathbf{m}|}{2} \int_0^s (\beta + S_u) du \right\} \frac{|\mathbf{m}|}{2} S_s \\
 &\quad \times \exp \left\{ -\frac{|\mathbf{m}| - 1}{2} \int_s^t (\beta + S_u) du \right\} ds \\
 &= \frac{|\mathbf{m}|}{2} \exp \left\{ -\frac{|\mathbf{m}|}{2} \int_0^t (\beta + S_u) du \right\} \\
 &\quad \times \int_0^t S_s \exp \left\{ \left(\frac{|\mathbf{m}|}{2} - \frac{|\mathbf{m}| - 1}{2} \right) \int_s^t (\beta + S_u) du \right\} ds \\
 &= |\mathbf{m}| \exp \left\{ -\frac{|\mathbf{m}|}{2} \int_0^t (\beta + S_u) du \right\} \\
 &\quad \times \left(1 - \exp \left\{ \left(\frac{|\mathbf{m}|}{2} - \frac{|\mathbf{m}| - 1}{2} \right) \int_0^t (\beta + S_u) du \right\} \right) \\
 &= |\mathbf{m}| \left(\exp \left\{ -\frac{|\mathbf{m}|}{2} \int_0^t (\beta + S_u) du \right\} \right. \\
 &\quad \left. - \exp \left\{ -\frac{|\mathbf{m}| - 1}{2} \int_0^t (\beta + S_u) du \right\} \right) \\
 &= |\mathbf{m}| \left(\frac{\beta}{(\beta + s)e^{\beta t/2} - s} \right)^{|\mathbf{m}| - 1} \left(1 - \frac{\beta}{(\beta + s)e^{\beta t/2} - s} \right).
 \end{aligned}$$

Iterating the argument leads to conclude that the death process jumps from $|\mathbf{m}|$ to $|\mathbf{m}| - i$ in $[0, t]$ with probability $\text{Bin}(|\mathbf{m}| - i \mid |\mathbf{m}|, p(t))$. \square

Note that when $s \in \mathbb{N}$, $\text{Ga}(\alpha_i + m, \beta + s)$ is the posterior of a $\text{Ga}(\alpha_i, \beta)$ prior given s Poisson observations with total count m . Hence the dual component $M_{i,t}$ is interpreted as the sum of the observed values of type i , and $S_t \subset \mathbb{R}_+$ as a continuous version of the sample size. In particular, (4.14) shows that a multivariate CIR propagates a vector of gamma distributions into a mixture whose kernels factorise into a gamma and a Dirichlet distribution, and whose mixing weights are driven by a one-dimensional death process with Binomial transitions together with hypergeometric probabilities for allocating the masses.

The following Proof of the conjugacy for mixtures of gamma random measures is due to Lo (1982) and outlined here for the ease of the reader.

Proof of Lemma 3.2. Since $z_{\mathbf{m}} := (z \mid \mathbf{m}) \sim \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta + s}$, from (2.6) we have

$$y_{m+1}, \dots, y_n \mid z, \mathbf{m}, n \stackrel{iid}{\sim} z_{\mathbf{m}} / |z_{\mathbf{m}}|, \quad n \mid z_{\mathbf{m}} \sim \text{Po}(|z_{\mathbf{m}}|).$$

Using (2.8) we have

$$\Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta + s} = \text{Ga}(\theta + |\mathbf{m}|, \beta + s) \Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}},$$

that is $|z_{\mathbf{m}}|$ and $z_{\mathbf{m}}/|z_{\mathbf{m}}|$ are independent with $\text{Ga}(\theta + |\mathbf{m}|, \beta + s)$ and $\Pi_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}$ distribution respectively. Then we have

$$z_{\mathbf{m}} \mid \mathbf{y}_{m+1:m+n} \sim \text{Ga}(\theta + |\mathbf{n}|, \beta + s + 1) \Pi_{\alpha + \sum_{i=1}^{K_{m+n}} n_i \delta_{y_i^*}} = \Gamma_{\alpha + \sum_{i=1}^{K_{m+n}} n_i \delta_{y_i^*}}^{\beta + s + 1}$$

where \mathbf{n} are the multiplicities of the distinct values in $\mathbf{y}_{1:n}$. Finally, by the independence of $|z_{\mathbf{m}}|$ and $z_{\mathbf{m}}/|z_{\mathbf{m}}|$, the conditional distribution of the mixing measure follows by the same argument used in Proposition 3.1. \square

We are now ready to prove the main result for DW processes.

Proof of Theorem 3.2. Fix a partition (A_1, \dots, A_K) of \mathcal{Y} . Then by Proposition 4.2

$$\begin{aligned} & \psi_t \left(\Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta + s} (A_1, \dots, A_K) \right) \\ &= \sum_{i=0}^{|\mathbf{m}|} \text{Bin}(|\mathbf{m}| - i; |\mathbf{m}|, p(t)) \text{Ga}(\theta + |\mathbf{m}| - i, \beta + S_t) \\ & \quad \times \sum_{\mathbf{0} \leq \mathbf{i} \leq \tilde{\mathbf{m}}, |\mathbf{i}|=i} p(\mathbf{i}; \tilde{\mathbf{m}}, i) \pi_{\alpha + \tilde{\mathbf{m}} - \mathbf{i}}, \end{aligned}$$

where $\Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta + s} (A_1, \dots, A_K)$ denotes $\Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta + s}(\cdot)$ evaluated on (A_1, \dots, A_K) and $\tilde{\mathbf{m}}$ are the multiplicities yielded by the projection of \mathbf{m} onto (A_1, \dots, A_K) . Use now (2.8) and (3.15) to write the right hand side of (3.14) as

$$\begin{aligned} & \sum_{\mathbf{n} \in L(\mathbf{m})} \tilde{p}_{\mathbf{m}, \mathbf{n}}(t) \Gamma_{\alpha + \sum_{i=1}^{K_m} n_i \delta_{y_i^*}}^{\beta + S_t} \\ &= \sum_{i=0}^{|\mathbf{m}|} \text{Bin}(|\mathbf{m}| - i; |\mathbf{m}|, p(t)) \text{Ga}(\theta + |\mathbf{m}| - i, \beta + S_t) \\ & \quad \times \sum_{\mathbf{0} \leq \mathbf{n} \leq \mathbf{m}, |\mathbf{n}|=i} p(\mathbf{n}; \mathbf{m}, i) \Pi_{\alpha + \sum_{j=1}^{K_m} (m_j - n_j) \delta_{y_j^*}}. \end{aligned}$$

Since the inner sum is the only term which depends on multiplicities and since Dirichlet processes are characterised by their finite-dimensional projections, we are only left to show that

$$\begin{aligned} & \sum_{\mathbf{0} \leq \mathbf{n} \leq \mathbf{m}, |\mathbf{n}|=i} p(\mathbf{n}; \mathbf{m}, i) \Pi_{\alpha + \sum_{j=1}^{K_m} (m_j - n_j) \delta_{y_j^*}} (A_1, \dots, A_K) \\ &= \sum_{\mathbf{0} \leq \mathbf{i} \leq \tilde{\mathbf{m}}, |\mathbf{i}|=i} p(\mathbf{i}; \tilde{\mathbf{m}}, i) \pi_{\alpha + \tilde{\mathbf{m}} - \mathbf{i}} \end{aligned}$$

which, in view of (4.7), holds if

$$\sum_{\mathbf{0} \leq \mathbf{n} \leq \mathbf{m}: \tilde{\mathbf{n}}=\mathbf{i}} p(\mathbf{i}; \mathbf{m}, i) = p(\mathbf{i}; \tilde{\mathbf{m}}, i),$$

where $\tilde{\mathbf{n}}$ denotes the projection of \mathbf{n} onto (A_1, \dots, A_K) . This is the consistency with respect to merging of classes of the multivariate hypergeometric distribution, and so the result now follows by the same argument at the end of the proof of Theorem 3.1. \square

We conclude by proving the recursive representation of Proposition 3.1, whose argument is analogous to the FV case.

Proof of Proposition 3.2. The update operation (3.17) follows directly from Lemma 3.1. The prediction operation (3.11) for elements of \mathcal{F}_Π follows from Theorem 3.2 together with the linearity of (1.4) and a rearrangement of the sums, so that

$$\begin{aligned} \psi_t & \left(\sum_{\mathbf{m} \in M} w_{\mathbf{m}} \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+s} \right) \\ &= \sum_{\mathbf{m} \in M} w_{\mathbf{m}} \sum_{\mathbf{n} \in L(\mathbf{m})} p_{\mathbf{m}, \mathbf{n}}(t) \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+S_t} \\ &= \sum_{\mathbf{n} \in L(M)} \left(\sum_{\mathbf{m} \in M, \mathbf{m} \geq \mathbf{n}} w_{\mathbf{m}} p_{\mathbf{m}, \mathbf{n}}(t) \right) \Gamma_{\alpha + \sum_{i=1}^{K_m} m_i \delta_{y_i^*}}^{\beta+S_t}. \end{aligned} \quad \square$$

As a final comment concerning the strategy followed for proving the propagation result in Theorems 3.1 and 3.2, one could be tempted to work directly with the duals of the FV and DW processes (Dawson and Hochberg, 1982; Ethier and Kurtz, 1993; Etheridge, 2000). However, this is not optimal, due to the high degree of generality of such dual processes. The simplest path for deriving the propagation step for the nonparametric signals appears to be resorting to the corresponding parametric dual by means of projections and by exploiting the filtering results for those cases.

Acknowledgements

The authors wish to thank an Associate Editor and two anonymous Referees for carefully reading the manuscript and for providing helpful comments.

References

- ANTONIAK, C. E. (1974). Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *Ann. Statist.* **2**, 1152–1174. [MR0365969](#)
- BARNDORFF-NIELSEN, O. and SHEPHARD, N. (2001). Non-Gaussian Ornstein-Uhlenbeck-based models and some of their uses in financial economics. *J. Roy. Statist. Soc. Ser. B* **63**, 167–241. [MR1841412](#)
- BEAL, M. J., GHAHRAMANI, Z. and RASMUSSEN, C. E. (2002). The infinite hidden Markov model. *Advances in Neural Information Processing Systems* **14**, 577–585.

- BLACKWELL, D. (1973). Discreteness of Ferguson selections. *Ann. Statist.* **1**, 356–358. [MR0348905](#)
- BLACKWELL, D. and MACQUEEN, J. B. (1973). Ferguson distributions via Pólya urn schemes. *Ann. Statist.* **1**, 353–355. [MR0362614](#)
- CARON, F., DAVY, M. and DOUCET, A. (2007). Generalized Pólya urn for time-varying Dirichlet process mixtures. *Proc. 23rd Conf. on Uncertainty in Artificial Intelligence*, Vancouver.
- CARON, F., DAVY, M., DOUCET, A., DUFLOS, E. and VANHEEGHE, P. (2008). Bayesian inference for linear dynamic models with Dirichlet process mixtures. *IEEE Trans. Sig. Proc.* **56**, 71–84. [MR2439814](#)
- CARON, F., NEISWANGER, W., WOOD, F., DOUCET, A. and DAVY, M. (2016). Generalized Pólya urn for time-varying Pitman–Yor processes. *J. Mach. Learn. Res.*, in press.
- CARON, F. and TEH, Y. W. (2012). Bayesian nonparametric models for ranked data. *Neural Information Processing Systems (NIPS 2012)*, Lake Tahoe, USA, 2012.
- COX, J. C., INGERSOLL, J. E. and ROSS, S. A. (1985). A theory of the term structure of interest rates. *Econometrica* **53**, 385–407. [MR0785475](#)
- CHALEYAT-MAUREL, M. and GENON-CATALOT, V. (2006). Computable infinite-dimensional filters with applications to discretized diffusion processes. *Stoch. Proc. Appl.* **116**, 1447–1467. [MR2260743](#)
- CHALEYAT-MAUREL, M. and GENON-CATALOT, V. (2009). Filtering the Wright–Fisher diffusion. *ESAIM Probab. Stat.* **13**, 197–217. [MR2518546](#)
- DALEY, D. J. and VERE-JONES (2008). *An introduction to the theory of point processes, Vol. 2*. Springer, New York.
- DAWSON, D. A. (1993). *Measure-valued Markov processes*. Ecole d’Eté de Probabilités de Saint Flour XXI. Lecture Notes in Mathematics **1541**. Springer, Berlin. [MR1242575](#)
- DAWSON, D. A. (2010). *Introductory lectures on stochastic population systems*. Technical Report Series **451**, Laboratory for Research in Statistics and Probability, Carleton University.
- DAWSON, D. A. and HOCHBERG, K. J. (1982). Wandering random measures in the Fleming–Viot model. *Ann. Probab.* **10**, 554–580. [MR0659528](#)
- DUNSON, D. B. (2006). Bayesian dynamic modeling of latent trait distributions. *Biostatistics* **7**, 551–568.
- ETHERIDGE, A. M. (2009). *Some mathematical models from population genetics*. École d’été de Probabilités de Saint-Flour XXXIX. Lecture Notes in Math. **2012**. Springer. [MR2759587](#)
- ETHERIDGE, A. M. (2000). *An introduction to superprocesses*. University Lecture Series, 20. American Mathematical Society, Providence, RI. [MR1779100](#)
- ETHIER, S. N. and GRIFFITHS, R. C. (1993). The transition function of a Fleming–Viot process. *Ann. Probab.* **21**, 1571–1590. [MR1235429](#)
- ETHIER, S. N. and GRIFFITHS, R. C. (1993b). The transition function of a measure-valued branching diffusion with immigration. In *Stochastic Processes. A Festschrift in Honour of Gopinath Kallianpur* (S. Cambanis, J. Ghosh, R. L. Karandikar and P. K. Sen, eds.), 71–79. Springer, New York. [MR1427302](#)

- ETHIER, S. N. and KURTZ, T. G. (1993). Fleming–Viot processes in population genetics. *SIAM J. Control Optim.* **31**, 345–386. [MR1205982](#)
- FAVARO, S., RUGGIERO, M. and WALKER, S. G. (2009). On a Gibbs sampler based random process in Bayesian nonparametrics. *Electron. J. Statist.* **3**, 1556–1566. [MR2578838](#)
- FERGUSON, T. S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209–230. [MR0350949](#)
- GASSIAT, E. and ROUSSEAU, J. (2016). Nonparametric finite translation hidden Markov models and extensions. *Bernoulli* **22**, 193–212. [MR3449780](#)
- GHOSAL, S. (2010). The Dirichlet process, related priors and posterior asymptotics. In *Bayesian Nonparametrics* (N. L. Hjort, C. C. Holmes, P. Müller and S. G. Walker, eds.). Cambridge Univ. Press, Cambridge [MR2730660](#)
- GRIFFIN, J. E. (2011). The Ornstein–Uhlenbeck Dirichlet Process and other time-varying processes for Bayesian nonparametric inference. *J. Stat. Plan. Inference.* **141**, 3648–3664. [MR2817371](#)
- GRIFFIN, J. E. and STEEL, M. F. J. (2006). Order-based dependent Dirichlet processes. *JASA* **473**, 179–194. [MR2268037](#)
- GRIFFITHS, R. C. and SPANÒ, D. (2010). Diffusion processes and coalescent trees. In *Probability and Mathematical Genetics: Papers in Honour of Sir John Kingman* (Bingham, N. H. and Goldie, C. M., eds.). London Mathematical Society Lecture Notes Series, Cambridge University Press. [MR2744247](#)
- GUTIERREZ, L., MENA, R. H. and RUGGIERO, M. (2016). A time dependent Bayesian nonparametric model for air quality analysis. *Comput. Statist. Data Anal.* **95**, 161–175. [MR3425946](#)
- JANSEN, S. and KURT, N. (2014). On the notion(s) of duality for Markov processes. *Probab. Surveys.* **11**, 59–120. [MR3201861](#)
- JOHNSON, N. L., KOTZ, S. and BALAKRISHNAN, N. (1997). *Discrete multivariate distributions*. John Wiley & Sons, New York. [MR1429617](#)
- KAWAZU, K. and WATANABE, S. (1971). Branching processes with immigration and related limit theorems. *Theory Probab. Appl.* **16**, 36–54. [MR0290475](#)
- KONNO, N. and SHIGA, T. (1988). Stochastic differential equations for some measure valued diffusions. *Probab. Th. Rel. Fields* **79**, 201–225. [MR0958288](#)
- LI, Z. (2011). *Measure-valued branching Markov processes*. Springer, Heidelberg. [MR2760602](#)
- LO, A. Y. (1982). Bayesian nonparametric statistical inference for Poisson point process. *Z. Wahrsch. Verw. Gebiete* **59**, 55–66. [MR0643788](#)
- MACEachern, S. N. (1999). Dependent Nonparametric Processes. In *ASA Proceedings of the Section on Bayesian Statistical Science*. American Statist. Assoc., Alexandria, VA.
- MACEachern, S. N. (2000). Dependent Dirichlet processes. *Tech. Rep.*, Ohio State University.
- MENA, R. H. and RUGGIERO, M. (2016). Dynamic density estimation with diffusive Dirichlet mixtures. *Bernoulli* **22**, 901–926. [MR3449803](#)
- MENA, R. H., RUGGIERO, M. and WALKER, S. G. (2011). Geometric stick-breaking processes for continuous-time Bayesian nonparametric modeling. *J. Statist. Plann. Inf.* **141**, 3217–3230. [MR2796026](#)

- PAPASPILIOPOULOS, O. and ROBERTS, G. O. (2008). Retrospective mcmc for dirichlet process hierarchical models. *Biometrika* **95**, 169–186. [MR2409721](#)
- PAPASPILIOPOULOS, O. and RUGGIERO, M. (2014). Optimal filtering and the dual process. *Bernoulli* **20**, 1999–2019. [MR3263096](#)
- RODRIGUEZ, A. and TER HORST, E. (2008). Bayesian dynamic density estimation. *Bayes. Anal.* **3**, 339–366. [MR2407430](#)
- RUGGIERO, M. and WALKER, S. G. (2009a). Bayesian nonparametric construction of the Fleming–Viot process with fertility selection. *Statist. Sinica* **19**, 707–720. [MR2514183](#)
- RUGGIERO, M. and WALKER, S. G. (2009b). Countable representation for infinite-dimensional diffusions derived from the two-parameter Poisson–Dirichlet process. *Elect. Comm. Probab.* **14**, 501–517. [MR2564485](#)
- STEPLETON, T., GHAHRAMANI, Z., GORDON, G., and LEE, T.-S. (2009). The block diagonal infinite hidden Markov model. *Journal of Machine Learning Research* **5**, 544–551.
- SETHURAMAN, J. (1994). A constructive definition of the Dirichlet process prior. *Statist. Sinica* **2**, 639–650. [MR1309433](#)
- SHIGA, T. (1990). A stochastic equation based on a Poisson system for a class of measure-valued diffusion processes. *J. Math. Kyoto Univ.* **30**, 245–279. [MR1068791](#)
- SPANÒ, D. and LIJOI, A. (2016). Canonical correlations for dependent gamma processes. [arXiv:1601.06079](#).
- TAVARÉ, S. (1984). Line-of-descent and genealogical processes, and their applications in population genetic models. *Theoret. Population Biol.* **26**, 119–164. [MR0770050](#)
- VAN GAEL, V., SAATCI, Y., TEH, Y. W. and GHAHRAMANI, Z. (2008). Beam sampling for the infinite hidden Markov model. In *Proceedings of the 25th International Conference on Machine Learning*.
- WALKER, S. G. (2007). Sampling the dirichlet mixture model with slices. *Comm. Statist. Sim. Comput.* **36**, 45–54. [MR2370888](#)
- WALKER, S. G., HATJISPYROS S. J. and NICOLERIS, T. (2007). A Fleming–Viot process and Bayesian nonparametrics. *Ann. Appl. Probab.* **17**, 67–80. [MR2292580](#)
- YAU, C., PAPASPILIOPOULOS, O., ROBERTS, G. O. and HOLMES, C. (2011). Bayesian non-parametric hidden Markov models with applications in genomics. *J. Roy. Statist. Soc. Ser. B* **73**, 37–57. [MR2797735](#)
- ZHANG, A., ZHU, J. and ZHANG, B. (2014). Max-margin infinite hidden Markov models. In *Proceedings of the 31st International Conference on Machine Learning*.