

A General results for Ornstein-Uhlenbeck processes

A.1 OU process with affine fitness

Proposition A.1. *Consider an OU process with affine fitness, whose density satisfies*

$$\partial_t p_t = -\nabla \cdot ((\mathbf{A}_t \mathbf{x} + \mathbf{e}_t) p_t - \mathbf{D}_t \nabla p_t) + (b_t + \mathbf{c}_t^\top \mathbf{x}) p_t. \quad (17)$$

where $\mathbf{A}_t \in \mathbb{R}^{d \times d}$, $\mathbf{e}_t, \mathbf{c}_t \in \mathbb{R}^d$, and $b_t \in \mathbb{R}$ are general, $\mathbf{D}_t \in \mathbb{R}^{d \times d}$ is symmetric positive definite. If $\rho_0 = m_0 \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$, then $\rho_t = m_t \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ for all $t \geq 0$, where

$$\dot{\boldsymbol{\Sigma}}_t = \mathbf{A}_t \boldsymbol{\Sigma}_t + \boldsymbol{\Sigma}_t \mathbf{A}_t^\top + 2\mathbf{D}_t, \quad (18)$$

$$\dot{\boldsymbol{\mu}}_t = \mathbf{A}_t \boldsymbol{\mu}_t + \boldsymbol{\Sigma}_t \mathbf{c}_t + \mathbf{e}_t, \quad (19)$$

$$\frac{\dot{m}_t}{m_t} = \mathbf{c}_t^\top \boldsymbol{\mu}_t + b_t. \quad (20)$$

Proof. We denote $\hat{p}_t(\mathbf{k})$ the Fourier transform of $p_t(\mathbf{x})$. We have the following identities

$$\mathcal{F}(\nabla \cdot (\mathbf{A}_t \mathbf{x} + \mathbf{e}_t) p_t) = -\mathbf{k}^\top \mathbf{A}_t \nabla_{\mathbf{k}} \hat{p}_t + i \mathbf{e}_t^\top \mathbf{k} \hat{p}_t, \quad (21)$$

$$\mathcal{F}(\nabla \cdot (\mathbf{D}_t \nabla p_t)) = -\mathbf{k}^\top \mathbf{D}_t \mathbf{k} \hat{p}_t. \quad (22)$$

Such that the Fourier transform of the PDE reads

$$\partial_t \hat{p}_t = (\mathbf{k}^\top \mathbf{A}_t + i \mathbf{c}_t^\top) \nabla_{\mathbf{k}} \hat{p}_t - \mathbf{k}^\top \mathbf{D}_t \mathbf{k} \hat{p}_t + (b_t - i \mathbf{e}_t^\top \mathbf{k}) \hat{p}_t. \quad (23)$$

We define the characteristics as the solution of

$$\frac{d\mathbf{k}_s}{ds} = -\mathbf{A}_s^\top \mathbf{k}_s - i \mathbf{c}_s, \quad (24)$$

such that

$$\frac{d\hat{p}_s(\mathbf{k}_s)}{ds} = \left(\frac{d\mathbf{k}_s}{ds} \right)^\top \nabla_{\mathbf{k}} \hat{p}_s(\mathbf{k}_s) + \partial_s \hat{p}_s(\mathbf{k}_s) \quad (25)$$

$$= -\mathbf{k}_s^\top \mathbf{D}_s \mathbf{k}_s \hat{p}_s + (b_s - i \mathbf{e}_s^\top \mathbf{k}_s) \hat{p}_s. \quad (26)$$

Denote $\Phi(t, t_0)$ the state transition matrix solution to the homogeneous ODE (24) with initial condition t_0 . It exists because the coefficient are locally integrable on \mathbb{R}^+ [41]. The state transition matrix has the following properties [42]:

$$\frac{d\Phi(s, u)}{ds} = -\mathbf{A}_s^\top \Phi(s, u), \quad \frac{d\Phi(s, u)}{du} = \Phi(s, u) \mathbf{A}_u^\top, \quad (27)$$

$$\Phi(t, u) \Phi(u, s) = \Phi(t, s), \quad \Phi(t, s)^{-1} = \Phi(s, t), \quad \Phi(s, s) = \mathbf{I}. \quad (28)$$

Denote $\Phi(s, 0) = \Phi_s$ for simplicity, we have, using the initial condition \mathbf{k}_0 :

$$\mathbf{k}_s = \Phi_s \mathbf{k}_0 - i \int_0^s \Phi(s, u) \mathbf{c}_u du. \quad (29)$$

We introduce $\tilde{\mathbf{v}}_s = -\int_0^s \Phi(s, u) \mathbf{c}_u du$, such that we have the Fourier transform along the characteristics reads, with an initial Gaussian distribution

$$\begin{aligned} \hat{p}_s(\mathbf{k}_s) = \exp \left[-i \mathbf{k}_0^\top \mathbf{x}_0 - \frac{1}{2} \mathbf{k}_0^\top \boldsymbol{\Sigma}_0 \mathbf{k}_0 - \mathbf{k}_0^\top \int_0^s \Phi_u^\top \mathbf{D}_u \Phi_u du \mathbf{k}_0 - 2i \mathbf{k}_0^\top \int_0^s \Phi_u^\top \mathbf{D}_u \tilde{\mathbf{v}}_u du \right. \\ \left. + \int_0^s \tilde{\mathbf{v}}_u^\top \mathbf{D}_u \tilde{\mathbf{v}}_u du + \int_0^s b_u du - i \int_0^s \mathbf{e}_u^\top \Phi_u \mathbf{k}_0 + \int_0^s \mathbf{e}_u^\top \tilde{\mathbf{v}}_u du \right]. \end{aligned} \quad (30)$$

We introduce the following quantities

$$\boldsymbol{\Psi}_s = \int_0^s \Phi_u^\top \mathbf{D}_u \Phi_u du, \quad \mathbf{q}_s^\top = \int_0^s \mathbf{e}_u^\top \Phi_u du, \quad \mathbf{u}_s = \int_0^s \Phi_u^\top \mathbf{D}_u \tilde{\mathbf{v}}_u du, \quad (31)$$

$$\gamma_s = \int_0^s \tilde{\mathbf{v}}_u^\top \mathbf{D}_u \tilde{\mathbf{v}}_u du, \quad w_s = \int_0^s \mathbf{e}_u^\top \tilde{\mathbf{v}}_u du. \quad (32)$$

Inverting the equation for characteristics we have $\mathbf{k}_0 = \Phi_s^{-1} \mathbf{k}_s + i \Phi_s^{-1} \int_0^s \Phi(s, u) \mathbf{c}_u du = \Phi_s^{-1} \mathbf{k}_s - i \Phi_s^{-1} \tilde{\mathbf{v}}_s$. We can directly see that the Fourier transform is quadratic in \mathbf{k} , such that it is the Fourier transform of a Gaussian measure. By considering one after the other the terms $O(k^2)$, $O(k)$ and $O(1)$, we can find analytical expressions for the covariance, the mean and the mass of the Gaussian measure as a function of the state transition matrix. They read:

$$\Sigma_s = \Phi_s^{-T} \Sigma_0 \Phi_s^{-1} + 2 \Phi_s^{-T} \Psi_s \Phi_s^{-1}, \quad (33)$$

$$\mu_s = \Phi_s^{-T} (\mu_0 + 2u_s + q_s) - \Sigma_s \tilde{\mathbf{v}}_s, \quad (34)$$

$$\log m_s = -\tilde{\mathbf{v}}_s^\top \mu_s - \frac{1}{2} \tilde{\mathbf{v}}_s^\top \Sigma_s \tilde{\mathbf{v}}_s + w_s + \gamma_s + \int_0^s b_u du. \quad (35)$$

Taking derivatives of these quantities and using the properties of the state transition matrix we recover the expected results. \square

A.2 OU process with quadratic fitness

Proposition 2.1 (OU process with quadratic fitness). *Consider an OU process with quadratic fitness, whose density satisfies*

$$\partial_t \rho_t(\mathbf{x}) = -\nabla \cdot ((\mathbf{A}_t \mathbf{x} + \mathbf{e}_t) \rho_t(\mathbf{x}) - \mathbf{D}_t \nabla \rho_t(\mathbf{x})) + (b_t + \mathbf{c}_t^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \Gamma_t \mathbf{x}) \rho_t(\mathbf{x}), \quad (8)$$

where $\mathbf{A}_t \in \mathbb{R}^{d \times d}$, $\mathbf{e}_t, \mathbf{c}_t \in \mathbb{R}^d$, and $b_t \in \mathbb{R}$ are generic, $\mathbf{D}_t \in \mathbb{R}^{d \times d}$ is symmetric positive definite, and $\Gamma_t \in \mathbb{R}^{d \times d}$ is symmetric negative semi-definite. If $\rho_0 = m_0 \mathcal{N}(\mu_0, \Sigma_0)$ with $m_0 > 0$, then $\rho_t = m_t \mathcal{N}(\mu_t, \Sigma_t)$ for all $t \geq 0$, where

$$\dot{\Sigma}_t = \mathbf{A}_t \Sigma_t + \Sigma_t \mathbf{A}_t^\top + 2\mathbf{D}_t + \Sigma_t \Gamma_t \Sigma_t, \quad (9)$$

$$\dot{\mu}_t = (\mathbf{A}_t + \Sigma_t \Gamma_t) \mu_t + \Sigma_t \mathbf{c}_t + \mathbf{e}_t, \quad (10)$$

$$\frac{\dot{m}_t}{m_t} = \frac{1}{2} \mu_t^\top \Gamma_t \mu_t + \mathbf{c}_t^\top \mu_t + b_t + \frac{1}{2} \text{tr}(\Gamma_t \Sigma_t). \quad (11)$$

Proof. Let's use the ansatz $p_t(\mathbf{x}) = u_t(\mathbf{x}) \exp[\mathbf{x}^\top \mathbf{G}_t \mathbf{x} / 2]$ where \mathbf{G}_t is symmetric. The equation for u_t now reads

$$\begin{aligned} \partial_t u_t + \frac{1}{2} u_t \mathbf{x}^\top \frac{d\mathbf{G}_t}{dt} \mathbf{x} = & -\nabla \cdot ((\mathbf{A}_t \mathbf{x} + \mathbf{e}_t - 2\mathbf{D}_t \mathbf{G}_t \mathbf{x}) u_t - \mathbf{D}_t \nabla u_t) \\ & + \left(b_t - \text{tr}(\mathbf{G}_t \mathbf{D}_t) + \mathbf{c}_t^\top \mathbf{x} - \mathbf{e}_t^\top \mathbf{G}_t \mathbf{x} + \mathbf{x}^\top \left(\frac{1}{2} \Gamma_t + \mathbf{G}_t \mathbf{D}_t \mathbf{G}_t - \mathbf{A}_t^\top \mathbf{G}_t \right) \mathbf{x} \right) u_t. \end{aligned} \quad (36)$$

The quadratic term vanishes when \mathbf{G}_t verifies for all \mathbf{x} the equation $\mathbf{x}^\top (\Gamma_t / 2 + \mathbf{G}_t \mathbf{D}_t \mathbf{G}_t - \mathbf{A}_t^\top \mathbf{G}_t) \mathbf{x} = \frac{1}{2} \mathbf{x}^\top \dot{\mathbf{G}}_t \mathbf{x}$. This is verified for \mathbf{G}_t satisfying the matrix Riccati equation

$$\Gamma_t + 2\mathbf{G}_t \mathbf{D}_t \mathbf{G}_t - \mathbf{A}_t^\top \mathbf{G}_t - \mathbf{G}_t \mathbf{A}_t = \frac{d\mathbf{G}_t}{dt}. \quad (37)$$

Provided the coefficients are locally integrable on \mathbb{R}^+ , that $\Gamma_t \leq 0$ and $\mathbf{D}_t \geq 0$, if $\mathbf{G}_0 \leq 0$, then there exists a unique solution \mathbf{G}_t on \mathbb{R}^+ for [37], and $\mathbf{G}_t \leq 0$ for all $t \geq 0$ [43]. We consider this unique solution with $\mathbf{G}_0 = 0$. Then, the equation for u_t becomes affine in growth

$$\partial_t u_t = -\nabla \cdot ((\tilde{\mathbf{A}}_t \mathbf{x} + \mathbf{e}_t - 2\mathbf{D}_t \mathbf{G}_t \mathbf{x}) u_t - \mathbf{D}_t \nabla u_t) + (b_t - \text{tr}(\mathbf{G}_t \mathbf{D}_t) + \tilde{\mathbf{c}}_t^\top \mathbf{x} - \mathbf{e}_t^\top \mathbf{G}_t \mathbf{x}) u_t. \quad (38)$$

We redefine $\tilde{\mathbf{A}}_t = \mathbf{A}_t - 2\mathbf{D}_t \mathbf{G}_t$, $\tilde{b}_t = b_t - \text{tr}(\mathbf{G}_t \mathbf{D}_t)$ and $\tilde{\mathbf{c}}_t = \mathbf{c}_t - \mathbf{G}_t \mathbf{e}_t$. The equation now reads

$$\partial_t u_t = -\nabla \cdot ((\tilde{\mathbf{A}}_t \mathbf{x} + \mathbf{e}_t) u_t - \mathbf{D}_t \nabla u_t) + (\tilde{b}_t + \tilde{\mathbf{c}}_t^\top \mathbf{x}) u_t. \quad (39)$$

We can apply the results of Prop. A.1. Since $u_0(\mathbf{x}) = p_0(\mathbf{x})$ is a Gaussian measure, it remains Gaussian, and we denote its covariance $\tilde{\Sigma}_t$, its mean $\tilde{\mathbf{x}}_t$ and its mass \tilde{m}_t . It follows that $p_t(\mathbf{x})$ is Gaussian measure for all $t \geq 0$ since $\tilde{\Sigma}_t^{-1} - \mathbf{G}_t$ is positive definite as the sum of positive definite and positive semi-definite terms. The covariance of this Gaussian measure is then $\Sigma_t = (\tilde{\Sigma}_t^{-1} - \mathbf{G}_t)^{-1}$.

We have the following fact for any differentiable matrix valued function B_t invertible for all t :

$$\frac{dB_t^{-1}}{dt} = -B_t^{-1} \frac{dB_t}{dt} B_t^{-1}. \quad (40)$$

Applying this to the covariance, we find:

$$\frac{d\Sigma_t}{dt} = -\Sigma_t \frac{d\Sigma_t^{-1}}{dt} \Sigma_t = \Sigma_t \tilde{\Sigma}_t^{-1} \left(\tilde{A}_t \tilde{\Sigma}_t + \tilde{\Sigma}_t \tilde{A}_t^\top + 2D_t \right) \tilde{\Sigma}_t^{-1} \Sigma_t + \Sigma_t \frac{dG_t}{dt} \Sigma_t. \quad (41)$$

Substituting $\tilde{A}_t = A_t - 2D_t G_t$ and $\dot{G}_t = \Gamma_t + 2G_t D_t G_t - A_t^\top G_t - G_t A_t$ we have

$$\frac{d\Sigma_t}{dt} = A_t \Sigma_t + \Sigma_t A_t^\top + 2D_t + \Sigma_t \Gamma_t \Sigma_t. \quad (42)$$

Completing the square in high-dimensions allows us to write the mean of the overall process as

$$\mu_t = \left(\tilde{\Sigma}_t^{-1} - G_t \right)^{-1} \tilde{\Sigma}_t^{-1} \tilde{\mu}_t = \Sigma_t \tilde{\Sigma}_t^{-1} \tilde{\mu}_t = (I + \Sigma_t G_t) \tilde{\mu}_t. \quad (43)$$

Using Prop. [A.1](#) we have

$$\frac{d\mathbf{x}_t}{dt} = \left(\Sigma_t G_t + \Sigma_t \frac{dG_t}{dt} \right) \tilde{\mu}_t + (I + \Sigma_t G_t) (\tilde{A}_t \tilde{\mu}_t + \tilde{\Sigma}_t \tilde{c}_t + \tilde{e}_t) \quad (44)$$

We replace $\tilde{A}_t = A_t - 2D_t G_t$ and we use the Riccati equation in the terms multiplying $\tilde{\mu}_t$. In the other terms we replace $\tilde{c}_t = c_t - G_t e_t$ and $e_t = \tilde{e}_t$. We find:

$$\frac{d\mu_t}{dt} = (A_t + \Sigma_t \Gamma_t) \mu_t + \Sigma_t c_t + e_t. \quad (45)$$

We can derive the fitness following the same approach and using lengthy simplifications. We only need to gather together the terms remaining after completing the square, as well as adjust for the change in covariance in the normalisation factor. However, this result is more easily obtained by using a result for the mean over a Gaussian probability distribution of quadratic form [\[44, Theorem 1.5\]](#). We find that

$$\frac{\dot{m}}{m} = \mathbb{E}[b_t + c_t^\top X + \frac{1}{2} X^\top \Gamma_t X] = b_t + c_t^\top \mu_t + \frac{1}{2} \text{tr}(\Gamma_t \Sigma_t) + \frac{1}{2} \mu_t^\top \Gamma_t \mu_t \quad (46)$$

where the expectation is understood to be taken with respect to $\mathcal{N}(x_t, \Sigma_t)$. \square

Corollary 2.2. *Consider an OU process with drift $v_t(x) = A_t x$ and a time-dependent fitness as in [\(8\)](#). Let $K_t \in \mathbb{R}^{n \times n}$ be an arbitrary matrix. Then there exists $\epsilon > 0$ such that the system is indistinguishable from another OU process with drift $v_t(x) = A_t + I + \epsilon K_t$ and time-dependent quadratic fitness.*

Proof. Let's define the following growth parameters: $\tilde{b}_t = b_t - \text{tr}((I + \epsilon K_t) \Sigma_t)/2$, $\tilde{c}_t = -(I + \epsilon K_t)^\top \Sigma_t^{-1} \mu_t$ and $\tilde{\Gamma}_t = \Gamma_t - ((I + \epsilon K_t)^\top \Sigma_t^{-1} + \Sigma_t^{-1} (I + \epsilon K_t))$. With these and the drift $\tilde{v}(x)$, the solution to the system of ODE above is unchanged. We take ϵ as the largest value such that $\tilde{\Gamma}_t$ is negative definite. This value is strictly larger than zero thanks to the identity. This derivation still holds if the drift is autonomous: if $v_t(x) = A x$, then there for any K , there exists $\epsilon > 0$ such that the system is indistinguishable from another OU process with $\tilde{v}_t(x) = A + I + \epsilon K$ and time-dependent quadratic fitness. \square

B Continuous-time loss for OU processes with quadratic fitness

Theorem 2.3 (Loss function for OU processes with quadratic fitness). *Consider the true process as well as the inferred process to both be OU processes with quadratic fitness, i.e. [\(8\)](#). Denoting*

$q = 2/\gamma$, with $K\Delta t$ fixed, when $\Delta t \rightarrow 0$ we have $\Delta t^{-1}L \rightarrow \mathcal{L}$, where \mathcal{L} is the continuous time loss function:

$$\mathcal{L} = q^{-1} \int_0^T dt m_t \left(\|\mathbf{v}_t - \hat{\mathbf{v}}_t\|_{\mathbf{X}_t^{-1}}^2 + \frac{1}{2} (h_t - \hat{h}_t)^2 \right) \quad (13)$$

$$+ q \sum_{i,j} \frac{\sigma_{i,t}^2 \sigma_{j,q,t}^2}{(\sigma_{j,q,t}^2 \sigma_{i,t}^2 + \sigma_{i,q,t}^2 \sigma_{j,t}^2)^2} \left(\mathbf{w}_{i,t}^\top (\mathbf{B}_t - \hat{\mathbf{B}}_t) \mathbf{w}_{j,t} \right)^2 + \lambda \int_0^T R_t dt \quad (14)$$

where $\Sigma_t = \sum_i \sigma_{i,t}^2 \mathbf{w}_{i,t} \mathbf{w}_{i,t}^\top$ is the eigendecomposition of the covariance Σ_t , $\sigma_{i,q,t}^2 = 1 + q\sigma_{i,t}^2$ for all i , $\mathbf{X}_t = 2\Sigma_t + q^{-1}$. We also have

$$\begin{aligned} \mathbf{B}_t - \hat{\mathbf{B}}_t &= \Sigma_t (\mathbf{A}_t - \hat{\mathbf{A}}_t)^\top + (\mathbf{A}_t - \hat{\mathbf{A}}_t) \Sigma_t + \Sigma_t (\Gamma_t - \hat{\Gamma}_t) \Sigma_t + 2(\mathbf{D}_t - \hat{\mathbf{D}}_t), \\ \mathbf{v}_t - \hat{\mathbf{v}}_t &= ((\mathbf{A}_t - \hat{\mathbf{A}}_t) + \Sigma_t (\Gamma_t - \hat{\Gamma}_t)) \boldsymbol{\mu}_t + \Sigma_t (\mathbf{c}_t - \hat{\mathbf{c}}_t) + (\mathbf{e}_t - \hat{\mathbf{e}}_t), \\ h_t - \hat{h}_t &= \frac{1}{2} \boldsymbol{\mu}_t^\top (\Gamma_t - \hat{\Gamma}_t) \boldsymbol{\mu}_t + (\mathbf{c}_t - \hat{\mathbf{c}}_t)^\top \boldsymbol{\mu}_t + (b_t - \hat{b}_t) + \frac{1}{2} \text{tr}((\Gamma_t - \hat{\Gamma}_t) \Sigma_t), \end{aligned} \quad (15)$$

and R_t is a strongly convex function of the parameters defining the fitness and the drift. For $\lambda > 0$, \mathcal{L} has a unique minimiser $t \mapsto (\hat{\mathbf{A}}_t^*, \hat{\mathbf{e}}_t^*, \hat{b}_t^*, \hat{\mathbf{c}}_t^*, \hat{\Gamma}_t^*)$.

Proof. When the entropic regularisation $\epsilon = 0$, the unbalanced Sinkhorn divergence between two Gaussian measures reduces to the Gaussian-Hellinger-Kantorovich distance $S_{0,\gamma} = \text{GHK}_\gamma$. Between the inferred and the true process at time t_i it reads:

$$\begin{aligned} &\text{GHK} \left(m_{t_i} \mathcal{N}(\boldsymbol{\mu}_{t_i}, \Sigma_{t_i}), \hat{m}_{t_i} \mathcal{N}(\hat{\boldsymbol{\mu}}_{t_i}, \hat{\Sigma}_{t_i}) \right) \\ &= 2q^{-1} \left(m_{t_i} + \hat{m}_{t_i} - 2 \sqrt{\frac{m_{t_i} \hat{m}_{t_i}}{\det \mathbf{J}} \exp \left[-\frac{1}{2} (\boldsymbol{\mu}_{t_i} - \hat{\boldsymbol{\mu}}_{t_i})^\top \mathbf{X}_{t_i}^{-1} (\boldsymbol{\mu}_{t_i} - \hat{\boldsymbol{\mu}}_{t_i}) \right]} \right), \end{aligned} \quad (47)$$

where we have

$$\mathbf{X}_{t_i} = \Sigma_{t_i} + \Sigma_{t_i} + q^{-1} \mathbf{I}, \quad \mathbf{J} = (\Sigma_{t_i,q} \Sigma_{t_i,q})^{1/2} (\mathbf{I} - q(\Sigma_{t_i} \Sigma_{t_i,q}^{-1} \Sigma_{t_i,q}^{-1} \Sigma_{t_i})^{1/2}), \quad (48)$$

with

$$\Sigma_{t_i,q} = q\Sigma_{t_i} + \mathbf{I}, \quad \Sigma_{t_i,q} = q\Sigma_{t_i} + \mathbf{I}. \quad (49)$$

Without loss of generality we consider the case $t_1 = \Delta t$, and we perform Taylor expansion in Δt of the GHK loss. Because the final expansion will be of order Δt^2 , only the terms of order Δt of the covariances play a role.

The hard part in this expansion is the expansion of $\det \mathbf{J}$. We denote $\Sigma_{\Delta t} = \Sigma_0 + \Delta t \mathbf{A}$ and $\tilde{\Sigma}_{\Delta t} = \Sigma_0 + \Delta t \tilde{\mathbf{A}}$ and $\Sigma_{0,q} = \mathbf{I} + q\Sigma_0$. For this, we need to compute $\Sigma_{\Delta t} \Sigma_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t} \tilde{\Sigma}_{\Delta t,q}^{-1}$, up to second order in Δt . We have that

$$\Sigma_{\Delta t} \Sigma_{\Delta t,q}^{-1} = (\Sigma_0 + \Delta t \mathbf{A}) (\Sigma_{0,q}^{-1} - q\Delta t \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1} + q^2 \Delta t^2 \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1}) \quad (50)$$

$$= \Sigma_0 \Sigma_{0,q}^{-1} + \Delta t (\mathbf{A} \Sigma_{0,q}^{-1} - q\Sigma_0 \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1}) \quad (51)$$

$$+ \Delta t^2 q (q\Sigma_0 \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1} - \mathbf{A} \Sigma_{0,q}^{-1} \mathbf{A} \Sigma_{0,q}^{-1}). \quad (52)$$

We also have

$$\tilde{\Sigma}_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t} = (\Sigma_{0,q}^{-1} - q\Delta t \Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} + q^2 \Delta t^2 \Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1}) (\Sigma_0 + \Delta t \tilde{\mathbf{A}}) \quad (53)$$

$$= \Sigma_{0,q}^{-1} \Sigma_0 + \Delta t (\Sigma_{0,q}^{-1} \tilde{\mathbf{A}} - q\Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} \Sigma_0) \quad (54)$$

$$+ \Delta t^2 q (q\Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} \Sigma_0 - \Sigma_{0,q}^{-1} \tilde{\mathbf{A}} \Sigma_{0,q}^{-1} \tilde{\mathbf{A}}). \quad (55)$$

We can now gather the $O(1)$, $O(\Delta t)$, $O(\Delta t^2)$ separately. We denote them respectively \mathbf{M} , \mathbf{H}_1 , \mathbf{H}_{2345} . We define the operation 'tt.' as the 'tilde transpose' operation, which is applied to

the term directly to its left. We therefore have

$$M = \Sigma_0 \Sigma_{0,q}^{-2} \Sigma_0 \quad (56)$$

$$H_1 = \Sigma_0 (\Sigma_{0,q}^{-2} \tilde{A} - q \Sigma_{0,q}^{-2} \tilde{A} \Sigma_{0,q}^{-1} \Sigma_0) + \text{tt}. \quad (57)$$

$$H_{2345} = q \Sigma_0 (q \Sigma_{0,q}^{-2} \tilde{A} \Sigma_{0,q}^{-1} \tilde{A} \Sigma_{0,q}^{-1} \Sigma_0 - \Sigma_{0,q}^{-2} \tilde{A} \Sigma_{0,q}^{-1} \tilde{A}) + \text{tt}. \quad (58)$$

$$+ A \Sigma_{0,q}^{-2} \tilde{A} + q^2 \Sigma_0 \Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-2} \tilde{A} \Sigma_{0,q}^{-1} \Sigma_0 \quad (59)$$

$$- q A \Sigma_{0,q}^{-2} \tilde{A} \Sigma_{0,q}^{-1} \Sigma_0 + \text{tt}. \quad (60)$$

Using the Woodbury formula we also have $\Sigma_0 \Sigma_{0,q}^{-1} = q^{-1}(\mathbf{I} - \Sigma_{0,q}^{-1})$. We have that Σ_0 and $\Sigma_{0,q}^{-1}$ commute, and that $M^{1/2} = \Sigma_0 \Sigma_{0,q}^{-1} = q^{-1}(\mathbf{I} - \Sigma_{0,q}^{-1})$. We introduce $S(Y)$ the unique solution X to the following Lyapunov equation

$$M^{1/2} X + X M^{1/2} = Y \quad (61)$$

This solution is expressed in terms of an integral (and is linear in Y), ie.

$$S(Y) = \int_0^\infty e^{-M^{1/2}t} Y e^{-M^{1/2}t} dt \quad (62)$$

As a result, we have

$$\left(\Sigma_{\Delta t} \Sigma_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t} \right)^{1/2} = M^{1/2} + \Delta t S(H_1) + \Delta t^2 (S(H_{2345}) - S(S(H_1)^2)) \quad (63)$$

We have $\mathbf{I} - q M^{1/2} = \Sigma_{0,q}^{-1}$, such that

$$\det \left(\mathbf{I} - q \left(\Sigma_{\Delta t} \Sigma_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t,q}^{-1} \tilde{\Sigma}_{\Delta t} \right)^{1/2} \right) = \det \Sigma_{0,q}^{-1} \det (\mathbf{I} + \Delta t \mathbf{L} + \Delta t^2 \mathbf{G}), \quad (64)$$

where $\mathbf{L} = -q \Sigma_{0,q} S(H_1)$ and $\mathbf{G} = -q \Sigma_{0,q} (S(H_{2345}) - S(S(H_1)^2))$. We then have the following expansion at second order in Δt

$$\det (\mathbf{I} + \Delta t \mathbf{L} + \Delta t^2 \mathbf{G}) = 1 + \Delta t \text{tr} \mathbf{L} + \frac{\Delta t^2}{2} (\text{tr}^2 \mathbf{L} - \text{tr} \mathbf{L}^2 + 2 \text{tr} \mathbf{G}). \quad (65)$$

So we need to compute $\text{tr} \mathbf{L}$, $\text{tr} \mathbf{G}$, $\text{tr} \mathbf{L}^2$. We have, using the fact that $\Sigma_{0,q}$ commutes with $M^{1/2}$,

$$\text{tr} \mathbf{L} = -q \text{tr} (\Sigma_{0,q} \int_0^\infty e^{-M^{1/2}t} H_1 e^{-M^{1/2}t} dt) = -\frac{q}{2} \text{tr} (M^{-1/2} \Sigma_{0,q} H_1) \quad (66)$$

$$= -\frac{q}{2} \text{tr} (\Sigma_{0,q}^2 \Sigma_0^{-1} H_1) = -\frac{q}{2} \left(\text{tr} (A + \tilde{A}) - q \text{tr} ((A + \tilde{A}) \Sigma_{0,q}^{-1} \Sigma_0) \right) \quad (67)$$

$$= -\frac{q}{2} \text{tr} (\Sigma_{0,q}^{-1} (A + \tilde{A})). \quad (68)$$

Computation of $\text{tr} \mathbf{G}$

Using the same trick as for $\text{tr} \mathbf{L}$ and standard properties of the trace we have

$$\text{tr} (-q \Sigma_{0,q} S(H_{2345})) = \frac{1}{2} \left(q^2 \text{tr} (A \Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} + \text{tt.}) - q \text{tr} (\Sigma_0^{-1} A \Sigma_{0,q}^{-2} \tilde{A}) \right) \quad (69)$$

To compute $\text{tr} (\Sigma_{0,q} S(S(H_1)^2))$ we denote $M^{1/2} = \sum_i u_i w_i w_i^\top$ the eigendecomposition of $M^{1/2}$. We also denote $\mathbf{W}_i = w_i w_i^\top$. Therefore we have

$$S(H_1) = \sum_{i,j} \frac{1}{u_i + u_j} \mathbf{W}_i H_1 \mathbf{W}_j \quad (70)$$

such that

$$\Sigma_{0,q} S(S(H_1)^2) = \sum_{i,j,k} \frac{(1 - u_k q)^2 (1 - u_j q)}{(u_i + u_k)(u_j + u_k)(u_i + u_j)} \mathbf{W}_i \left(u_i \tilde{A} + u_k A \right) \mathbf{W}_k \left(u_k \tilde{A} + u_j A \right) \mathbf{W}_j. \quad (71)$$

Taking the trace we have

$$\text{tr}(\Sigma_{0,q} S(S(H_1)^2)) = \sum_{i,j} \frac{(1-u_j q)^2 (1-u_i q)}{2u_i(u_i + u_j)^2} \text{tr} \left(W_i(u_i \tilde{A} + u_j A) W_j(u_j \tilde{A} + u_i A) \right). \quad (72)$$

Additionally, we have that

$$\text{tr} \left(W_i(u_i \tilde{A} + u_j A) W_j(u_j \tilde{A} + u_i A) \right) \quad (73)$$

$$= u_i u_j \text{tr} \left(W_i(A - \tilde{A}) W_j(A - \tilde{A}) \right) + (u_i + u_j)^2 \text{tr} \left(W_i A W_j \tilde{A} \right) \quad (74)$$

As a result we have

$$\text{tr} G = \frac{1}{2} q^2 \text{tr}(A \Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} + \text{tt.}) + \frac{q}{2} \sum_{i,j} \frac{(1-u_j q)^2 (1-u_i q) u_j}{(u_i + u_j)^2} \text{tr} \left(W_i(A - \tilde{A}) W_j(A - \tilde{A}) \right) \quad (75)$$

Computation of $\text{tr} L^2$

Similarly we have

$$\Sigma_{0,q} S(H_1) = \sum_{i,j} \frac{(1-u_j q)}{(u_i + u_j)} W_i \left(u_i \tilde{A} + u_j A \right) W_j. \quad (76)$$

Using the same approach for $\text{tr} G$ we find

$$\text{tr}((\Sigma_{0,q} S(H_1))^2) = \sum_{i,j} \frac{(1-u_i q)(1-u_j q) u_i u_j}{(u_i + u_j)^2} \text{tr} \left(W_i(A - \tilde{A}) W_j(A - \tilde{A}) \right) \quad (77)$$

$$+ \text{tr}(\Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} \tilde{A}). \quad (78)$$

As a result we have

$$\text{tr} L^2 = q^2 \text{tr}(\Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} \tilde{A}) + q^2 \sum_{i,j} \frac{(1-u_i q)(1-u_j q) u_i u_j}{(u_i + u_j)^2} \text{tr} \left(W_i(A - \tilde{A}) W_j(A - \tilde{A}) \right) \quad (79)$$

Computation of $\det(\Sigma_{\Delta t,q} \tilde{\Sigma}_{\Delta t,q})^{1/2}$

We have

$$\Sigma_{\Delta t,q} \tilde{\Sigma}_{\Delta t,q} = \Sigma_{0,q}^2 \left(I + q \Delta t (\Sigma_{0,q}^{-2} A \Sigma_{0,q} + \Sigma_{0,q}^{-1} \tilde{A}) + q^2 \Delta t^2 \Sigma_{0,q}^{-2} A \tilde{A} \right). \quad (80)$$

Taking the determinant and Taylor expanding we have

$$\det(\Sigma_{\Delta t,q} \tilde{\Sigma}_{\Delta t,q}) = \det(\Sigma_{0,q})^2 \quad (81)$$

$$\times \left(1 + q \Delta t \text{tr}(\Sigma_{0,q}^{-1} (A + \tilde{A})) + q^2 \frac{\Delta t^2}{2} \left((\text{tr}^2(\Sigma_{0,q}^{-1} (A + \tilde{A})) - \text{tr}(\Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} A + \text{tt.})) \right) \right). \quad (82)$$

Expanding the square root we find

$$\det(\Sigma_{\Delta t,q} \tilde{\Sigma}_{\Delta t,q})^{1/2} = \det(\Sigma_{0,q}) \quad (83)$$

$$\times \left(1 + q \frac{\Delta t}{2} \text{tr}(\Sigma_{0,q}^{-1} (A + \tilde{A})) + q^2 \frac{\Delta t^2}{8} \left(\text{tr}^2(\Sigma_{0,q}^{-1} (A + \tilde{A})) - 2 \text{tr}(\Sigma_{0,q}^{-1} A \Sigma_{0,q}^{-1} A + \text{tt.}) \right) \right) \quad (84)$$

Computation of $\det \mathbf{J}$

Going back to $\det \mathbf{J}$, we are left with

$$\det \mathbf{J} = \left(1 + q \frac{\Delta t}{2} \text{tr}(\Sigma_{0,q}^{-1}(\mathbf{A} + \tilde{\mathbf{A}})) + q^2 \frac{\Delta t^2}{8} \left(\text{tr}^2(\Sigma_{0,q}^{-1}(\mathbf{A} + \tilde{\mathbf{A}})) - 2\text{tr}(\Sigma_{0,q}^{-1}\mathbf{A}\Sigma_{0,q}^{-1}\mathbf{A} + \text{tt.}) \right) \right) \quad (85)$$

$$\times \left(1 + \Delta t \text{tr} \mathbf{L} + \frac{\Delta t^2}{2} (\text{tr}^2 \mathbf{L} - \text{tr} \mathbf{L}^2 + 2\text{tr} \mathbf{G}) \right). \quad (86)$$

Using the result for $\text{tr} \mathbf{L}$ we see that the terms in Δt cancel, and the final expansion is of order Δt^2 . As a result, after simplifications we are left with

$$\det \mathbf{J} = \left(1 - q^2 \frac{\Delta t^2}{4} \text{tr}(\Sigma_{0,q}^{-1}\mathbf{A}\Sigma_{0,q}^{-1}\mathbf{A} + \text{tt.}) + \frac{\Delta t^2}{2} (-\text{tr} \mathbf{L}^2 + 2\text{tr} \mathbf{G}) \right). \quad (87)$$

Simplifying further we find that

$$\det \mathbf{J} = 1 + \frac{\Delta t^2 q}{2} \left(\frac{q}{2} \text{tr}(\Sigma_{0,q}^{-1}(\mathbf{A} - \tilde{\mathbf{A}})\Sigma_{0,q}^{-1}(\mathbf{A} - \tilde{\mathbf{A}})) \right) \quad (88)$$

$$+ \sum_{i,j} \frac{u_j(1-u_iq)(1-u_jq)(1-q(u_j+u_i))}{(u_i+u_j)^2} \text{tr} \left(\mathbf{W}_i(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{W}_j(\mathbf{A} - \tilde{\mathbf{A}}) \right) \quad (89)$$

This can be simplified

$$\left(\frac{q}{2} \text{tr}(\Sigma_{0,q}^{-1}(\mathbf{A} - \tilde{\mathbf{A}})\Sigma_{0,q}^{-1}(\mathbf{A} - \tilde{\mathbf{A}})) \right) \quad (90)$$

$$+ \sum_{i,j} \frac{u_j(1-u_iq)(1-u_jq)(1-q(u_j+u_i))}{(u_i+u_j)^2} \text{tr} \left(\mathbf{W}_i(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{W}_j(\mathbf{A} - \tilde{\mathbf{A}}) \right) \quad (91)$$

$$= \sum_{i,j} \frac{(1-qu_i)(1-qu_j)}{(u_i+u_j)^2} \left(\frac{q}{2} u_i^2 - \frac{q}{2} u_j^2 + u_j \right) \text{tr} \left(\mathbf{W}_i(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{W}_j(\mathbf{A} - \tilde{\mathbf{A}}) \right). \quad (92)$$

which leaves us with the following simplification

$$\det \mathbf{J} = 1 + \frac{\Delta t^2 q}{2} \sum_{i,j} \frac{u_i(1-u_iq)(1-u_jq)}{(u_i+u_j)^2} \text{tr} \left(\mathbf{W}_i(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{W}_j(\mathbf{A} - \tilde{\mathbf{A}}) \right). \quad (93)$$

Using $\Sigma_0 = \sum_i \sigma_i^2 \mathbf{W}_i$, we can compute the eigenvalue decomposition as a function of σ_i . Using the notation $\sigma_{i,q}^2 = 1 + q\sigma_i^2$ we have

$$\frac{u_i(1-u_iq)(1-u_jq)}{(u_i+u_j)^2} = \frac{\sigma_i^2 \sigma_{j,q}^2}{(\sigma_{j,q}^2 \sigma_i^2 + \sigma_{i,q}^2 \sigma_j^2)^2}. \quad (94)$$

Finally, at first non-zero order in Δt the determinant reads

$$\det \mathbf{J} = 1 + \frac{\Delta t^2 q}{2} \sum_{i,j} \frac{\sigma_i^2 \sigma_{j,q}^2}{(\sigma_{j,q}^2 \sigma_i^2 + \sigma_{i,q}^2 \sigma_j^2)^2} \text{tr} \left(\mathbf{W}_i(\mathbf{A} - \tilde{\mathbf{A}})\mathbf{W}_j(\mathbf{A} - \tilde{\mathbf{A}}) \right). \quad (95)$$

Taylor expansion of the GHK term

We now expand the masses and mean at first order in Δt , $m_{\Delta t} = m_0 + \Delta t m_0 h$, $\tilde{m}_{\Delta t} = m_0 + \Delta t m_0 \tilde{h}$, $\boldsymbol{\mu}_{\Delta t} = \boldsymbol{\mu}_0 + \Delta t \mathbf{v}$, $\tilde{\boldsymbol{\mu}}_{\Delta t} = \boldsymbol{\mu}_0 + \Delta t \tilde{\mathbf{v}}$. We have at first non zero order in Δt

$$\exp \left[-\frac{1}{2} (\boldsymbol{\mu}_{\Delta t} - \tilde{\boldsymbol{\mu}}_{\Delta t})^\top \mathbf{X}_{\Delta t}^{-1} (\boldsymbol{\mu}_{\Delta t} - \tilde{\boldsymbol{\mu}}_{\Delta t}) \right] = 1 - \frac{\Delta t^2}{2} (\mathbf{v} - \tilde{\mathbf{v}})^\top \mathbf{X}_0^{-1} (\mathbf{v} - \tilde{\mathbf{v}}). \quad (96)$$

We denote $\mathbf{X}_0 = 2\mathbf{\Sigma}_0 + q^{-1}$. Expanding the remaining terms, the zeroth and first order terms cancel, leading to the expansion

$$\text{GHK} \left(m_{\Delta t} \mathcal{N}(\boldsymbol{\mu}_{\Delta t}, \mathbf{\Sigma}_{\Delta t}), \hat{m}_{\Delta t} \mathcal{N}(\hat{\boldsymbol{\mu}}_{\Delta t}, \hat{\mathbf{\Sigma}}_{\Delta t}) \right) \quad (97)$$

$$= m_0 q^{-1} \Delta t^2 \left(\left((\mathbf{v} - \tilde{\mathbf{v}})^\top \mathbf{X}_0^{-1} (\mathbf{v} - \tilde{\mathbf{v}}) + \frac{1}{2} (h - \tilde{h})^2 \right) \right) \quad (98)$$

$$+ q \sum_{i,j} \frac{\sigma_i^2 \sigma_{j,q}^2}{(\sigma_{j,q}^2 \sigma_i^2 + \sigma_{i,q}^2 \sigma_j^2)^2} \text{tr} \left(\mathbf{W}_i (\mathbf{A} - \tilde{\mathbf{A}}) \mathbf{W}_j (\mathbf{A} - \tilde{\mathbf{A}}) \right). \quad (99)$$

Integrating over all snapshots, the continuous time loss reads

$$\begin{aligned} \mathcal{L} = & \int_0^T m_t q^{-1} \left(\left((\mathbf{v} - \tilde{\mathbf{v}})^\top \mathbf{X}_0^{-1} (\mathbf{v} - \tilde{\mathbf{v}}) + \frac{1}{2} (h - \tilde{h})^2 \right) \right. \\ & + q \sum_{i,j} \frac{\sigma_i^2 \sigma_{j,q}^2}{(\sigma_{j,q}^2 \sigma_i^2 + \sigma_{i,q}^2 \sigma_j^2)^2} \text{tr} \left(\mathbf{W}_i (\mathbf{A} - \tilde{\mathbf{A}}) \mathbf{W}_j (\mathbf{A} - \tilde{\mathbf{A}}) \right) \left. \right) dt \\ & + \lambda \int_0^T \int \rho_t(\mathbf{x}) (\|\mathbf{v}_t(\mathbf{x})\|_2^2 + \alpha |g_t(\mathbf{x})|^2) d\mathbf{x} dt. \end{aligned} \quad (100)$$

The first order expansions \mathbf{A} , h , and \mathbf{v} are obtained directly using the ODEs in Prop. 2.1, giving the final result.

Let's denote $\theta_t \in \mathbb{R}^N$ the vector of the N parameters defining the growth and the drift at time t . We study the functional $\mathcal{L}[\theta_t]$ which reads

$$\mathcal{L}[\theta_t] = \int_0^T (F_t(\theta_t) + \lambda R_t(\theta_t)) dt. \quad (101)$$

where F is the part of the integrand coming from the expansion of the Gaussian-Hellinger-Kantorovich, and R is the regularisation.

Let's take $t \in [0, T]$. Because $\theta \mapsto F_t(\theta)$ is the composition of convex functions and of affine maps, $\theta \mapsto F_t(\theta)$ is also convex. Let's show that $\theta \mapsto R_t(\theta)$ is a strongly convex function.

Lemma B.1. *Let $w : (\theta, x) \in \mathbb{R}^N \times \mathbb{R}^n \mapsto \mathbb{R}^d$ be a function continuous in x and linear in θ . Let $f : \mathbb{R}^d \mapsto \mathbb{R}$ be a continuous, strictly convex function and ρ a continuous function from \mathbb{R}^n to $]0, +\infty[$. We then have that*

$$\ell : \theta \in \mathbb{R}^N \mapsto \int \rho(x) f(w(\theta, x)) dx \quad (102)$$

is strictly convex.

Proof. Let's take $\theta_1 \neq \theta_2$ and $\alpha \in]0, 1[$. By linearity, we have $\forall x$:

$$w(\alpha\theta_1 + (1-\alpha)\theta_2, x) = \alpha w(\theta_1, x) + (1-\alpha)w(\theta_2, x). \quad (103)$$

Because $\theta \mapsto f(w(\theta, x))$ is the composition of a convex function and of a linear map, it is convex for all x . Then, $\forall x$,

$$f(w(\alpha\theta_1 + (1-\alpha)\theta_2, x)) = f(\alpha w(\theta_1, x) + (1-\alpha)w(\theta_2, x)) \leq \alpha f(w(\theta_1, x)) + (1-\alpha)f(w(\theta_2, x)). \quad (104)$$

Because w is linear in θ , each of its component in \mathbb{R}^d is a multivariate polynomial function of θ of degree one. By uniqueness of the coefficient of polynomial functions, because $\theta_1 \neq \theta_2$ there exists x_0 such that $w(\theta_1, x_0) \neq w(\theta_2, x_0)$. By continuity, this is also true on an open set $U \subset \mathbb{R}^n$ centred in x_0 . Therefore, since f is strictly convex, the inequality is strict for all $x \in U$, i.e.:

$$f(w(\alpha\theta_1 + (1-\alpha)\theta_2, x)) = f(\alpha w(\theta_1, x) + (1-\alpha)w(\theta_2, x)) < \alpha f(w(\theta_1, x)) + (1-\alpha)f(w(\theta_2, x)). \quad (105)$$

By multiplying by $\rho(x)$, which is strictly positive, and by integrating, we keep the inequality strict and we find

$$\ell(\alpha\theta_1 + (1-\alpha)\theta_2) < \alpha\ell(\theta_1) + (1-\alpha)\ell(\theta_2), \quad (106)$$

proving that ℓ is strictly convex. \square

	Hidden layers	Batch size	Learning rate	Iterations
Bistable	64, 64, 64	256	0.01	25,000
Bifurcating CLE	64, 64, 64	256	0.01	10,000
Haematopoietic CLE	128, 128, 128	256	0.01	10,000
Lineage tracing	128, 128, 128	256	0.01	10,000

Table 5: Hyperparameter settings: score networks

		Hidden (force)	Hidden (growth)	Batch	LR	Iterations	λ	α	γ
UPFI	Bistable	64, 64, 64	64, 64, 64	256	0.003	5,000	0.01	0.1	5
	Bif CLE	64, 64, 64	64	256	0.003	10,000	0.001	1	5
	Haem CLE	128, 128, 128	128	256	0.003	10,000	0.001	1	5
	Lineage	128, 128, 128	128, 128, 128	256	0.001	10,000	0.001	1	25
PFI	Bistable	64, 64, 64	–	256	0.003	5,000	0.01	–	–
	Bif CLE	64, 64, 64	–	256	0.003	10,000	0.001	–	–
	Haem CLE	128, 128, 128	–	256	0.003	10,000	0.001	–	–
	Lineage	128, 128, 128	–	256	0.001	10,000	0.001	–	–
ODE	Bistable	Coupled to growth	64, 64, 64	256	0.003	5,000	0.01	0.1	5
	Lineage	Coupled to growth	128, 128, 128	256	0.001	10,000	0.001	0.1	25
TIGON++	Bistable	64, 64, 64	64, 64, 64	256	0.003	5,000	0.01	0.1	5
	Lineage tracing	128, 128, 128	128, 128, 128	256	0.001	10,000	0.001	0.1	25

Table 6: Hyperparameter settings: dynamics

We have $\rho_t(x) > 0$ for all x because it is a Gaussian density. The drift and the fitness being linear in θ , this shows that

$$\theta \mapsto R_t(\theta) = \int \rho_t(x) (\|v_t(x)\|_2^2 + \alpha |g_t(x)|^2) dx \quad (107)$$

is a strictly convex function of θ . Additionally, it is a multivariate polynomial function of θ of degree two, so its Hessian is a definite positive and constant, proving that $\theta \mapsto R_t(\theta)$ is strongly convex, and so is $F_t + R_t$. Along with the fact that both $(t, \theta) \mapsto F_t(\theta)$ and $(t, \theta) \mapsto R_t(\theta)$ are continuous in $[0, T] \times \mathbb{R}^N$, this is enough to ensure the existence and uniqueness of a minimum $t \mapsto \theta_t^*$ for the functional $\mathcal{L}[\theta_t]$ [45]. \square

C Implementation and experiment details

We implement Alg. 1 using PyTorch and employ the GeomLoss package [46] for computation of the unbalanced Sinkhorn divergence. All model training was carried out using a NVIDIA L40S GPU.

C.1 Score matching

While in principle any score matching approach to learn $s_t(x) = \nabla \log p_t(x)$ can be used within Alg. 1, in practice we employ denoising score matching [24] within the noise-conditional score network framework introduced in [5]. For $K + 1$ snapshots taken at times $(t_i)_{i=0}^K$, we parameterise the time-dependent score using a multilayer perceptron (MLP) $s_\phi(t, x, \eta) = \text{NN}_\phi(d + 2, d)$ where d is the dimension and η is the noise level, and train using the algorithm described in [5, Section 4.2]. While a range of noise levels $\eta_0 < \dots < \eta_L$ are used for training the score, subsequently for training the probability flow we use the smallest noise scale $\eta = \eta_0$, representing our final estimate of the score.

In what follows, we use a default of $L = 5$ noise levels logarithmically spaced between $(\exp(-2), 1)$. Score networks $s_t(x)$ are parameterised using MLPs with ReLU activations. In all case, score training was carried out using noise-conditional denoising score matching [5, 24] using the AdamW optimiser with the hyperparameter choices listed in Table 5.

C.2 Training: UPFI and PFI

Additive noise models For UPFI, in all cases we parameterise an *autonomous* force $v_\theta(x)$ using a MLP. This is because, in all simulated systems, the true force is also autonomous. In the experimental lineage tracing dataset we reason that an autonomous force would be consistent with the biologically motivated model of a Waddington’s landscape [29]. We parameterise separately the force $v_\theta(x)$

and growth $g_\theta(\mathbf{x})$ using MLPs with ReLU activations. We train UPFI models using the AdamW optimiser, our architecture and hyperparameter choices are listed in Table 6.

For PFI, motivated by the observations in the Gaussian case we reason that an autonomous force is insufficient to fit the data if growth is not accounted for. We therefore employ a *non-autonomous* force, parameterising $\mathbf{v}_\theta(t, \mathbf{x}) = \text{NN}_\theta(d+1, d)(t, \mathbf{x})$ with ReLU activations. We train PFI models using the AdamW optimiser, our architecture and hyperparameter choices are listed in Table 6.

Multiplicative noise model For the multiplicative noise models we considered in Fig. 5, we parameterise an autonomous force in two components, corresponding to production and degradation terms in the model (16). Specifically, we let

$$\mathbf{f}_\theta(\mathbf{x}) = \text{NN}_\theta(d, d)(\mathbf{x}), \quad \mathbf{g}_\theta(\mathbf{x}) = \text{NN}_\theta(d, d)(\mathbf{x}),$$

and the resulting force is $(\mathbf{f}_\theta - \mathbf{g}_\theta)(\mathbf{x})$. To constrain the output of both networks to be non-negative, we opt for a Softplus activation on the final layer of outputs. For all other layers we use the ReLU activation as a default choice. Architectures for (\mathbf{f}, \mathbf{g}) and all other hyperparameter choices are as given in Table 6.

C.3 Training: TIGON++

The problem of dynamical transport for systems with mass imbalance was previously studied in the work [16]. In this work, the authors consider *deterministic* systems only and allow both the force $\mathbf{v}_t(\mathbf{x})$ and growth $g_t(\mathbf{x})$ to be time dependent. However, the TIGON algorithm relies on kernel density estimation (KDE) from the input data [16, Methods] which is sensitive to the choice of kernel bandwidth and suffers from the curse of dimensionality as the dimension increases. The choice of data-fitting loss, in the form of minimising squared discrepancies between the predicted and KDE densities, adds to these difficulties: this loss relies pointwise on estimated densities and is thus not “geometry-aware” in the sense of optimal transport based losses [46]. Finally, propagating densities under flow models e.g. (4) are well known to be computationally costly. For these reasons, training the TIGON algorithm was infeasible for most of our numerical experiments.

Because we needed a robust comparison baseline, we decided to implement the same model used by TIGON (deterministic transport with growth), but train it with the UPFI training procedure. With UPFI, we circumvent the need for density estimation by using the probability flow formulation of the Fokker-Planck equation. Together with the use of the unbalanced Sinkhorn divergence as the data fitting loss, we believe that this substantially improves the TIGON method and also makes for a more rigorous baseline to test against UPFI. We call TIGON++ this re-implementation of TIGON.

Specifically, we parameterise a drift $\mathbf{v}_\theta(t, \mathbf{x}) = \text{NN}_\theta(d+1, d)(t, \mathbf{x})$ and growth $g_\theta(t, \mathbf{x}) = \text{NN}_\theta(d+1, 1)(t, \mathbf{x})$ with ReLU activations. For a sampled data point $(\mathbf{x}_0, m_0 = 1)$ at $t = 0$, its state $(\mathbf{x}_{t_i}, m_{t_i})$ at each timepoint t_i is simulated by forward integration of the system

$$\dot{\mathbf{x}}_t = \mathbf{v}_\theta(t, \mathbf{x}_t), \quad \dot{m}_t = g_\theta(t, \mathbf{x}_t)m_t, \quad (108)$$

and we form the empirical distribution $\hat{\rho}_{t_i}(\mathbf{x}) = \sum_{k=1}^{N_i} m_{k, t_i} \delta(\hat{\mathbf{x}}_{k, t_i} - \mathbf{x})$ for each timepoint t_i . For the rest of the training procedure we use the same data-fitting loss as UPFI, i.e. (12).

C.4 Training: fitness-ODE

Motivated by issues pertaining to the identifiability of dynamics involving both drift and growth, we propose a well-known dynamical model as a baseline model for inference. Let $\{\rho_t\}_t$ be a continuous distributional path satisfying mild conditions in the space of measures describing some population evolution. Then there exists a unique scalar field $U_t(\mathbf{x}_t)$ such that ρ_t satisfies

$$\partial_t \rho_t(\mathbf{x}) = -\nabla \cdot (\rho_t(\mathbf{x}) \nabla U_t(\mathbf{x})) + U_t(\mathbf{x}) \rho_t(\mathbf{x}). \quad (109)$$

That this is the case can be read from [47, Section A.3] or [48, Proposition 2.2]. This can be interpreted as a continuous dynamics where $U_t(\mathbf{x})$ is the *fitness* of state \mathbf{x} . The rate at which agents reproduce is prescribed by $U_t(\mathbf{x})$, and agents migrate to regions of higher fitness following $\nabla U_t(\mathbf{x})$. Theoretically, for a regular enough sequence of population snapshots $\{\rho_t\}_t$, a single time-dependent fitness function U_t is sufficient to generate the path $t \mapsto \rho_t$ via these dynamics. While it is perhaps not obvious, a *single* quantity, the fitness U_t , is enough to generate the full path ρ_t in the space

of measures. As a baseline, we therefore propose a neural parameterisation of U_t and to learn U_t following the TIGON++ setup, but with $v_t(\mathbf{x}) = \nabla U_t(\mathbf{x})$ and $g_t(\mathbf{x}) = U_t(\mathbf{x})$. We parameterise $U_t(\mathbf{x}) = \text{NN}_\theta(d+1, 1)(t, \mathbf{x})$ using ReLU activations. All hyperparameter choices are listed in Table 6.

C.5 DeepRUOT

We use the existing DeepRUOT implementation provided by [32], which parameterises the force, growth rate and score function. Different to our method, however, the individual components of the dynamics are coupled via a physics-informed neural network (PINN)-type loss ([32, Section 5.3]) that aims to incorporate information from the governing Fokker-Planck equation. For the bistable system example (Fig. 3) we use their PyTorch implementation of [32, Algorithm 1]. For each of the drift, growth and score networks, three hidden layers of size 128 were used. For further details on DeepRUOT implementation and training we refer the reader to [32] and accompanying code.

C.6 Forward simulation

Bistable system We consider a potential-driven dynamics in dimension $d \in \{2, 5, 10\}$ specified by

$$\begin{aligned} \mathbf{v} &= -\nabla V, & V(\mathbf{x}) &= 0.9\|\mathbf{x} - \mathbf{a}\|_2^2\|\mathbf{x} - \mathbf{b}\|_2^2 + 10\sum_{i=3}^d x_i^2 \\ b(\mathbf{x}) &= \frac{5}{2}(1 + \tanh(2x_0)), \\ d(\mathbf{x}) &= 0. \end{aligned}$$

For $t \in [0, 1]$ we simulate particles following $d\mathbf{X}_t = \mathbf{v}(\mathbf{X}_t) dt + \sigma d\mathbf{B}_t$ using the Euler-Maruyama method. We set the noise level to $\sigma = 1/2$, and use the initial condition $\mathbf{X}_0 \sim \mathcal{N}(\mathbf{0}, 0.01\mathbf{I})$. At each Euler step, simulated particles divide with probability $b(\mathbf{X})\Delta t$. We simulate starting from $N_0 = 500$ particles, and the total population size grows over time following the prescribed dynamics. Population snapshots are taken from *independent* realisations of the process at $K+1 = 5$ timepoints uniformly spaced between $[0, 1]$.

Reaction network systems The bifurcating and HSC reaction networks were taken from previous literature [34, 35], corresponding to the networks BF and HSC in the collection of BoolODE benchmarking problems [34]. The original implementation, however, modelled both gene and protein expression levels and as a result does not strictly fall in the modelling framework we consider. This is because protein levels are not observed and thus are hidden variables. We re-implemented each of these systems to involve only gene expression dynamics, and also change the noise model: in the original implementation an ad-hoc square-root noise model was used, i.e. $\sigma(\mathbf{x}) = \alpha\sqrt{\mathbf{x}}$. We choose to use a more biophysically motivated noise model ([16]), and take

$$\sigma_i(\mathbf{x}) = \sqrt{f_i(\mathbf{x}) + \lambda_i x_i},$$

where $\mathbf{f}(\mathbf{x})$ is a vector-valued function of state-dependent production rates for each gene x_i , and each λ_i is the corresponding degradation rate. For the growth rates, we consider a scenario where cells in one branch of the system trajectory divide at a faster rate than the others:

- In the bifurcating network, we set

$$\beta(x) = 5 \left(\frac{\tanh(5(x_7 - 0.7)) + 1}{2} \right) \left(1 - \frac{\tanh(5(x_1 - 0.7)) + 1}{2} \right)$$

- In the HSC network, we set

$$\beta(x) = 5.5 \left(\frac{\tanh(x_E - 1.5) + 1}{2} \right).$$

We simulate both systems using the same code as for the bistable system, with the volume parameter $1/\sqrt{V} = 0.5$, starting with a population of 500 cells and capturing $K+1 = 10$ timepoints. For bifurcating and HSC network we use simulation time intervals of $0 \leq t \leq 1.25$ and $0 \leq t \leq 1$ respectively.

Fate probability computation We provide a straightforward definition of fate probability in what follows. Let \mathbf{x}_t be the state of an observed cell or individual at time t . Let $\sqcup_i \Omega_i$ be a partitioning of the state space (e.g. some subset of \mathbb{R}^d) where each Ω_i is understood to correspond to well-defined, stable states of the system at some final time, say $t = 1$. In the biological setting, this is typically thought of as a mature cell “type” [1, 29]. Then the fate probability of \mathbf{x}_t towards Ω_i is defined as the conditional probability:

$$\mathbb{P}(\mathbf{X}_1 \in \Omega_i | \mathbf{X}_t = \mathbf{x}_t).$$

In practice, such as for the bistable system of Fig. 3, we form a partitioning of the state space into two regions by running k -means with $k = 2$ on the final snapshot from the system. Given any query state \mathbf{x}_t at time t , we empirically estimated true and inferred fate probabilities by forward simulation of either the ground truth SDE or inferred dynamics:

$$\mathbb{P}(\mathbf{X}_1 \in \Omega_i | \mathbf{X}_t = \mathbf{x}_t) \approx M^{-1} \sum_{i=1}^M \mathbf{1}_{\Omega_i}(\mathbf{X}_1 | \mathbf{X}_t = \mathbf{x}_t),$$

where M is the number of trials to sample.

C.7 Neural graphical model

For the Neural Graphical Model (NGM) example of Fig. 4, we use the architecture introduced in [36] as a drop-in parameterisation of the autonomous force $\mathbf{v}_\theta(\mathbf{x})$. For each output variable, we use two hidden layers with sizes [64, 64]. We use a group lasso regularisation strength $\lambda_{\text{GL}} = 0.03$ and employ the proximal update scheme outlined in [36] Section C.4.1] with a learning rate of 0.003 and train for 5,000 iterations. We use the score networks that were already pre-trained for the additive and multiplicative UPFI models. All other training details are taken to be the same as for the earlier UPFI training.

C.8 Single cell lineage tracing data

Preprocessing Data for the study of [1] are available from the original publication using the GEO database with accession number GSE140802. Starting from raw counts, expression data is normalised using `dyn.pp.recipe_monocle` function from the Dynamo package [37]. In brief, raw gene expression values are per-cell normalised and then $\log(1 + x)$ -transformed. For all our experiments we use the 10-dimensional PCA embedding of cell gene expression profiles. Spliced and unspliced transcript counts were obtained from reanalysis of the raw sequencing data of [1] and RNA velocity estimates were subsequently obtained using the Dynamo package [37]. Scripts and datasets for this re-analysis are available upon request.

From the full dataset, 86,416 cells deemed to be contributing to the “Neutrophil-Monocyte” trajectory (as determined by the original publication [1]) were selected. Using the 10-dimensional PCA embedding for these data, we apply UPFI, PFI, fitness-ODE and TIGON. We do not include DeepRUOT in this analysis since it resulted in an out-of-memory error in the initial stages of training. Noting also that the original publication [32] considered only the 2D SPRING layout, we believe that further modification of its training pipeline may be necessary.

For UPFI, we train a time-dependent score model with hidden dimensions [128, 128, 128] for 10,000 iterations with a batch size of 256 and learning rate of 10^{-2} . We adopt an additive noise model and parameterise an autonomous force $\mathbf{v}_\theta(\mathbf{x})$ and growth $g_\theta(\mathbf{x})$ each with a MLP with hidden dimensions [128, 128, 128]. We set $\gamma = 25.0$, $\lambda = 0.001$, $\alpha = 1.0$ and we set $\sigma = 0.5$. Note that the choice of γ is not scale-invariant, and we found that the typical length scale in the lineage tracing data is larger than in the simulation data. We train for 10,000 iterations with a batch size of 256 and learning rate 10^{-3} .

We train PFI with the same hyperparameter choices as UPFI, except we use a non-autonomous force as done in earlier examples. Finally, we train TIGON and fitness-ODE following the training procedure outlined in Sections C.3 and C.4 and the same hyperparameters as for UPFI and PFI.

C.9 Remark on UPFI in the case when only frequencies are available

When N_i/N_0 is not a good estimator for the ratio $|\rho_{t_i}|/|\rho_0|$, we can only build an estimator for the normalised density $\rho_{t_i}/|\rho_{t_i}|$:

$$\frac{\rho_{t_i}}{|\rho_{t_i}|} \simeq \frac{1}{N_i} \sum_{k=1}^{N_i} \delta(\mathbf{x} - \mathbf{x}_{k,t_i}) \quad (110)$$

This poses limitations on the fitness which can be inferred. Indeed, writing $\tilde{\rho}_t = \rho_t/|\rho_t|$, we have

$$\partial_t \tilde{\rho}_t = \partial_t \left(\frac{\rho_t}{|\rho_t|} \right) = \frac{1}{|\rho_t|} \partial_t \rho_t - (\partial_t \log |\rho_t|) \tilde{\rho}_t.$$

Substituting back the PDE governing ρ_t , we find that $\tilde{\rho}_t$ is also governed by a drift-diffusion PDE, but with a time-dependent bias in the source term:

$$\partial_t \tilde{\rho}_t = -\nabla \cdot [\tilde{\rho}_t (\mathbf{v}_t - \nabla \cdot \mathbf{D}_t - \mathbf{D}_t \nabla \log \tilde{\rho}_t)] + (g_t - \partial_t \log |\rho_t|) \tilde{\rho}_t.$$

Therefore, when we don't have access to the absolute number of individuals present in a population at a given time, we can only hope to infer the fitness g_t up to a time-dependent bias. In this case, the UPFI approach can also be applied using a large enough mass conservation strength q in the unbalanced Sinkhorn distance.