# An unsupervised machine-learning checkpoint-restart algorithm using Gaussian mixtures for particle-in-cell simulations

G. Chen*, L. Chacón, T. B. Nguyen

*Los Alamos National Laboratory, Los Alamos, NM 87545*

## Abstract

We propose an unsupervised machine-learning checkpoint-restart (CR) algorithm for particle-in-cell (PIC) algorithms using Gaussian mixtures (GM). The algorithm features a particle compression stage and a particle reconstruction stage, where a continuum particle distribution function (PDF) is constructed and resampled, respectively. To guarantee fidelity of the CR process, we ensure the exact preservation of invariants such as charge, momentum, and energy for both compression and reconstruction stages, everywhere on the mesh. We also ensure the preservation of Gauss' law after particle reconstruction. As a result, the GM CR algorithm is shown to provide a clean, conservative restart capability while potentially affording orders of magnitude savings in input/output requirements. We demonstrate the algorithm using a recently developed exactly energy- and charge-conserving PIC algorithm using both electrostatic and electromagnetic tests. The tests demonstrate not only a high-fidelity CR capability, but also its potential for enhancing the fidelity of the PIC solution for a given particle resolution.

*Keywords:* unsupervised machine learning, Gaussian mixture model, particle-in-cell, checkpoint restart,
*PACS:*

## 1. Introduction

Resiliency, data locality, and asynchrony are key major challenges facing the practical use of exascale computing for scientific applications. Because of extreme concurrency, very large system scale, and complex memory hierarchies, hardware failures (both "soft" and "hard") are expected to become more frequent towards and beyond exascale. Currently, 100 billion transistors/node , thousands of nodes, 10M-core supercomputers are built (e.g. Summit and Sierra [1]). The very large total number of components will lead to frequent failures, even though the mean time between failures (MTBF) for the individual components may be large. For instance, while the MTBF of a CPU can be months to years [2], that of current supercomputers can be within a few hours [3, 4]. With billion-core parallelism at exascale, the MTBF has been projected to be within (or even far below) one hour [5, 6]. Therefore, it is important to enable efficient strategies that allow software and algorithms to perform in a frequently interrupted environment.

---

*Corresponding author
   *Email address:* gchen@lanl.gov (G. Chen)

Particle-based simulation algorithms are widely employed, at the heart of many algorithmic strategies (e.g., Monte Carlo, particle-in-cell, molecular dynamics) and applications (e.g., aerosol transport in combustion and climate, radiation transport, and plasma transport). Checkpoint/restart (CR) enables simulation recovery from previous interrupted simulations due to either finite queue wall-clock-time limits or hardware (HW) failures. This is commonly done by storing a sufficiently complete data snapshot to disk at given time intervals, which can be then read back to restart the simulation. Particle-based simulations at the extreme scale are particularly challenged by the input/output (IO) requirements of storing billions to trillions of particles, as is already the case in the leading-class supercomputers. The challenges are significantly worsened by the current trend towards hierarchical architectures, featuring many levels of parallelism, each delivered by different architectural solutions. Synchronous checkpointing in hierarchical systems would require bulk synchronization across the levels of the hierarchy, and ultimately storage in the file system via IO. Asynchronous IO, as well as memory-based IO, are being explored as partial solutions to the CR problem [7]. Nevertheless, *any* IO-based CR strategy would greatly benefit from a high-fidelity compression strategy of data for particle simulations.

In this study, we explore the viability of an *unsupervised machine-learning, optimization-based CR strategy* for plasma particle-in-cell (PIC) simulations, combining *optimal* (in some sense, to be clarified below) compression and reconstruction of particle data. Compression of particle data is performed per spatial cell by construction of a continuum particle distribution function (PDF) with a Gaussian mixture [8], based on a penalized maximum-likelihood-estimation (PMLE) approach using complexity criteria [9]. The resulting optimization problem is solved by an adaptive Expectation-Maximization (EM) algorithm [10, 9], which can automatically search for the optimal number of Gaussian components satisfying a generalized "minimum-message-length (MML)" Bayesian Information Criterion [11]. The method can be formulated to conserve up to second moments exactly [12]. Particle-data is reconstructed (also locally per cell) by sampling of the PDF in velocity space (here using Monte Carlo), with a simple moment-matching projection technique [13]. Particle spatial positions within a given cell are re-initialized randomly (i.e, we assume that the plasma is uniform within a cell). Both compression and reconstruction operations are local in configuration space (i.e., each computational cell features an *independent* PDF reconstruction process) and done *in-situ* (assuming that the cell has sufficient particles, e.g., more than 10), and only Gaussian parameters are checkpointed.

A quiescent restart in plasmas requires, in addition to the preservation of (at least) moments up to second order, the enforcement of Gauss' law (i.e., $\nabla \cdot \mathbf{E} = \rho$ where $\mathbf{E}$ is electric field, and $\rho$ is charge density) discretely everywhere on the spatial mesh. Gauss' law is closely related to charge conservation, and local violations will result in plasma waves being launched to equilibrate charge. The electric field is saved at CR, and is thus available at both compression and reconstruction stages. To enforce Gauss' law discretely, it is sufficient to ensure resampled particles *exactly* match the charge density field per species and per cell. We accomplish this by correcting particle weights according to a straightforward mass-matrix solve [14].

The potential for IO compression of particle data using GM is quite large. Each Gaussian component of the mixture requires ten parameters to be fully specified, which is comparable to the number of degrees of freedom needed per particle in a 3D-3V PIC method (e.g., three positions and three velocities, plus particle weight and optionally an integer identifying the cell on the mesh). Given that a few Gaussians (< 10) are usually sufficient (as demonstrated in our numerical tests) to capture most details of the PDF, and that typical PIC simulations employ hundreds if not thousands of particles per cell, it follows that GM can easily result in several orders of magnitude savings in IO requirements for checkpointing particle data.

The proposed CR strategy exactly conserves local (per cell) charge, momentum, and energy, satisfies Gauss' law everywhere, is massively parallel, communication-avoiding, locality-aware, and asynchronous by construction (except for the mass-matrix solve step), and only synchronizes and checkpoints compressed data. It is worth pointing out that we are not the first ones to realize the potential of GM PDF reconstruction in PIC algorithms, with various authors having used it in the past for diagnostics [15], to couple with other physical processes [16], or, more related to this study, for Gaussian-to-Gaussian remapping in 1D-1V phase-space to eliminate Gaussian-shape distortion in a finite-mass-method-based Vlasov-Poisson algorithm [17]. However, to our knowledge, this is the first application of an *adaptive* GM algorithm for particle data compression in CR of PIC simulations.

The rest of the paper is organized as follows. Section 2 introduces the basic concepts of the PMLE method employed in this study to learn and resample the Gaussian mixture, including the strategies to enforce Gauss' law and conserve of up to second moments. Section 3 demonstrates the CR algorithm for prototypical plasma-physics electrostatic and electromagnetic PIC tests, and demonstrate the potential of GM to improve the PIC solution for a given particle resolution. We also explore ways to improve the efficiency of the underlying EM algorithm to find the GM. Finally, we conclude in Section 4.

## 2. Methodology

We describe the two main elements of the CR GM strategy, namely, GM component estimation (particle-data compression) and GM sampling (particle-data reconstruction). Specifically, an adaptive EM algorithm is used to estimate the number of components of the Gaussian mixture and their parameters, and a moment-matching sampling technique is used to regenerate particles from the Gaussian mixture.

A GM is defined as a convex combination of *K* Gaussian distributions:

$$f(\mathbf{x}) = \sum_{k=1}^{K} \omega_k f_k(\mathbf{x}), \tag{1}$$

where each Gaussian $f_k$ is weighted by $\omega_k$ with $\sum_k \omega_k = 1$ and $w_k > 0$. The Gaussian distribution is defined as

$$f_k(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^D |\overline{\overline{\mathbf{\Sigma}}}_k|}} e^{-(\mathbf{x}-\mathbf{\mu}_k)^T \overline{\overline{\mathbf{\Sigma}}}_k^{-1}(\mathbf{x}-\mathbf{\mu}_k)/2}, \tag{2}$$

where $\mathbf{\mu}$ is a $D$-dimensional mean vector, $\overline{\overline{\mathbf{\Sigma}}}$ is a $D \times D$ covariance matrix, and $|\overline{\overline{\mathbf{\Sigma}}}|$ is the determinant of $\overline{\overline{\mathbf{\Sigma}}}$.

3

## 2.1. Adaptive GM models and the penalized maximum likelihood function

The goal is to estimate the parameters $\theta \equiv \{\omega, \mu, \overset{\leftrightarrow}{\Sigma}\}$ of each Gaussian as well as the number of mixture components, $K$, given $N$ independent samples $\mathbf{X} = (\mathbf{x}_1...\mathbf{x}_N)$ drawn from $f(\mathbf{x})$. Conventionally, maximum likelihood is used to estimate $\theta$ for a prescribed number of components [18]. However, estimating the number of components itself is in fact also important, and can be addressed in the framework of the Bayesian Information Criterion [8]. In what follows, we give a brief overview of this approach.

We seek to find the maximum likelihood of the model $K$ (given a data set $\mathbf{X}$), which by Bayes' rule reads:

$$p(K|\mathbf{X}) = \frac{p(\mathbf{X}|K)p(K)}{p(\mathbf{X})}, \tag{3}$$

where $p(K)$ is the prior probability distribution for the model family, and $p(\mathbf{X}) = \sum p(\mathbf{X}|K)p(K)$ is a normalizing constant. If we assume that all models are equally likely *a priori*, then $p(K)$ is uniform. Therefore, maximizing $p(K|\mathbf{X})$ is equivalent to maximizing $p(\mathbf{X}|K)$, which is the so-called marginal likelihood (also known as evidence [19] or type II maximum likelihood [20]), and is given by:

$$p(\mathbf{X}|K) = \int p(\mathbf{X}|\theta, K)p(\theta|K)d\theta, \tag{4}$$

where $p(\theta|K)$ is a prior probability distribution, and $p(\mathbf{X}|\theta, K)$ is the likelihood function, which for a Gaussian mixture reads:

$$p(\mathbf{X}|\theta, K) = \sum_{k=1}^{K} \omega_k f_k(\mathbf{x}_i|\mu_k, \overline{\overline{\Sigma}}_k).$$

We seek to maximize Eq. 4. For completeness, the derivation is carried out in Appendix A, and results in the penalized log-likelihood function:

$$L(\theta) = \ln\left[p(\mathbf{X}|\theta, K)\right] - \frac{d}{2}\ln N - \frac{T}{2}\sum_{i=1}^{K}\ln(\omega_i), \tag{5}$$

The last term of Eq. 5 is crucial for finding the number of components, and avoiding over-fitting and singularities of standard maximum likelihood estimate (MLE) [9]. As pointed out in Ref. [9], this term is an effective Dirichlet prior with negative parameters. Such a prior has a strong tendency to annihilate redundant components. We refer to Ref. [21] for a theoretical treatment on this important point, and Refs. [22, 23] for its practical use in the context of Gaussian mixtures.

## 2.2. Learning the GM model by a penalized MLE

We follow the standard method of MLE to seek optimum values of the Gaussian parameters. This is achieved by maximizing the penalized likelihood function, Eq. 5. For a given set of particles, the penalized log-likelihood function is given by

$$L(\theta) = \sum_{p=1}^{N} \alpha_p \ln\left[\sum_{k=1}^{K} \omega_k f_k(\mathbf{v}_p|\mu_k, \overline{\overline{\Sigma}}_k)\right] - \frac{d}{2}\ln N - \frac{T}{2}\sum_{k=1}^{K}\ln(\omega_k), \tag{6}$$

where $\mathbf{v}_p$ is particle velocity and $\alpha_p$ is the particle weight, which accounts for cases with non-identical samples [24]. Note that $\sum_{p=1}^{N} \alpha_p = N$. Typically, the MLE estimator is found by solving the likelihood equation:

$$\frac{\partial \ln L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0, \tag{7}$$

subject to the constraint that $\sum_k \omega_k = 1$, with:

$$\frac{\partial^2 \ln L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^2} < 0. \tag{8}$$

Setting the derivative of Eq. 6 with respect to the mean $\boldsymbol{\mu}_k$ of the Gaussian components to zero, we obtain

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{p=1}^{N} \gamma_{pk} \mathbf{v}_p, \tag{9}$$

where

$$\gamma_{pk} \equiv \frac{\alpha_p \omega_k f_k(\mathbf{v}_p | \boldsymbol{\mu}_k, \bar{\bar{\boldsymbol{\Sigma}}}_k)}{\sum_{k=1}^{K} \omega_k f_k(\mathbf{v}_p | \boldsymbol{\mu}_k, \bar{\bar{\boldsymbol{\Sigma}}}_k)}, \tag{10}$$

and

$$N_k = \sum_{p=1}^{N} \gamma_{pk}. \tag{}$$

Note that

$$\sum_{k=1}^{K} \gamma_{pk} = \alpha_p. \tag{11}$$

Setting the derivative of Eq. 6 with respect to $\boldsymbol{\Sigma}_k$ to zero, we obtain

$$\bar{\bar{\boldsymbol{\Sigma}}}_k = \frac{1}{N_k} \sum_{p=1}^{N} \gamma_{pk} (\mathbf{v}_p - \boldsymbol{\mu}_k)(\mathbf{v}_p - \boldsymbol{\mu}_k)^{\mathrm{T}}. \tag{12}$$

Maximizing Eq. 6 with respect to the mixing coefficients (again, constrained by $\sum_k \omega_k = 1$) gives [25, 9]:

$$\omega_k = \frac{N_k - \frac{T}{2}}{N - \frac{T}{2}K}, \tag{13}$$

provided $N_k - \frac{T}{2} > 0$. This suggests one should begin with more components than the "true" number of components of the mixture [21]. A component is eliminated ($K^{new} \leftarrow K^{old} - 1$) if $N_k - \frac{T}{2} \leq 0$. In the limit of $N \to \infty$, Eq. 13 recovers the standard MLE result, i.e.,

$$\tilde{\omega}_k = \frac{N_k}{N}. \tag{14}$$

The solution to Eqs. 9-13 can only be found iteratively.

*2.3. GM component estimation: Expectation-Maximization algorithm (EM-GM)*

The EM-GM algorithm provides an iterative procedure to find a local maximum of the log-likelihood function with respect to the Gaussian parameters. For an extensive review of theoretical and practical aspects of EM algorithm for finite mixture models, see Ref [26]. As discussed above, we start with a relatively

large number of Gaussians (~10), for each Gaussian, we set its mean to coincide with a randomly chosen particle, and its variance to be the same as the total variance. Each EM-GM iteration consists of the following steps [27]:

1. For each Gaussian component $k$, perform E-step: Given the parameter set $\theta_k^{it}$, where the superscript $it$ denotes the iteration level, evaluate Eq. 10.

2. For the same Gaussian component $k$, perform M-step: Compute $\theta_k^{it+1} = \{\omega, \mu, \overline{\overline{\Sigma}}\}_k^{it+1}$ via Eqs 9,12, and 13. If $\omega_k \leq 0$, remove the Gaussian, and let $K^{it+1} = K^{it} - 1$, otherwise, $K^{it+1} = K^{it}$.

3. Re-normalize weight by $\omega_k = \omega_k / \sum_{i=1}^{K^{it+1}} \omega_i$.

4. Repeat steps 1 to 3 until all Gaussians are updated, check for convergence by monitoring the log-likelihood function, Eq. 6.

Depending on how much the overlap of the Gaussians, convergence of the algorithm may be slow. In general, convergence is slow when Gaussians are poorly separated, and one should consider ways to accelerate it for practical applications (see results in Sec. 3).

### 2.3.1. Properties of the EM-GM algorithm

An important property of the EM-GM algorithm based on the *unpenalized* MLE is that it conserves up to second moments of the sample particles *exactly*, i.e., the mass, mean, and variance of the mixture coincide with those of sample particles (see Ref. [12] for the 1D MLE, and derivation below for the multivariate case). As a consequence, physical quantities such as mass, momentum, and energy (or more precisely the pressure tensor) are conserved by the GM continuum reconstruction of the particle PDF. However, such conservation property is not inherited by the EM-GM algorithm based on the *penalized* MLE (for component adaptivity) . To recover the moment conservation property, which is desirable for high-fidelity CR in particle simulations, we perform the estimate in two steps: first use the PMLE to select the optimal number of Gaussians, and then postprocess the result with one step of unpenalized MLE to regain conservation.

It is useful to derive the conservation properties of the unpenalized-MLE-based EM-GM algorithm as follows. We begin with the conservation of the first moment (mean):

$$
\begin{aligned}
E(\mathbf{v}) &= \sum_{k=1}^{K} \omega_k \boldsymbol{\mu}_k \\
&= \sum_{k=1}^{K} \omega_k \frac{1}{N_k} \sum_{p=1}^{N} \gamma_{pk} \mathbf{v}_p \\
&= \frac{1}{N} \sum_{p=1}^{N} \alpha_p \mathbf{v}_p = E(\mathbf{v}_p),
\end{aligned}
\tag{15}
$$

where we use the law of total expectation [28] for the first equality, Eq. 9 for the second equality, and Eq. 14 and 11 for the third equality. It is easily seen that, if Eq. 13 is used instead of 14, the third equality above would not follow through, breaking conservation of the first moment.

6

The derivation of the preservation of the second moments (variance) follows a similar procedure,

$$
\begin{aligned}
Var(\mathbf{v}) &= E(Var(\mathbf{v}|\mathbf{y})) + Var(E(\mathbf{v}|\mathbf{y})) \\
&= \sum_{k=1}^{K} \omega_k \Sigma_k + E(E(\mathbf{v}|\mathbf{y})^2) - E(E(\mathbf{v}|\mathbf{y}))^2 \\
&= \sum_{k=1}^{K} \omega_k \frac{1}{N_k} \sum_{p=1}^{N} \gamma_{pk}(\mathbf{v}_p - \boldsymbol{\mu}_k)(\mathbf{v}_p - \boldsymbol{\mu}_k)^{\mathrm{T}} + \sum_{k=1}^{K} \omega_k \mu_k^2 - E(x)^2 \\
&= \frac{1}{N} \sum_{p=1}^{N} \alpha_p v_p^2 - \left( \frac{1}{N} \sum_{p=1}^{N} \alpha_p \mathbf{v}_p \right)^2 = Var(\mathbf{v}_p),
\end{aligned}
$$

where $\mathbf{y}$ is the hidden variable indicating the Gaussian component that a particle belongs to. Here, the first equality is the law of total variance [28], the second equality uses definitions of expectations and variances, and the third equality uses Eq. 12 and the so-called Adam's law [i.e., $E(\mathbf{v}) = E(E(\mathbf{v}|\mathbf{y}))$] [28]. To get to the fourth equality, Eqs. 14, 11, and 15 are used. We observe that using Eq. 13 instead of Eq. 14 would again break the equality of the variance between the Gaussians and particles.

The derivations above indicate that Eq. 14 is critical for the conservation properties we wish to preserve for the Gaussians. Equation 13 is, however, critical for selecting the correct number of Gaussian components. We have designed a procedure that combines the advantages of both (i.e., conservation *and* adaptivity) as follows. We first perform iterations using Eq. 13 to prune out unnecessary Gaussians. Once converged, we simply perform an extra step using Eq. 14 for the mixing coefficients. This is equivalent to accounting for the Gaussian weights based only on the data, without penalization. Once we have Eq. 14 satisfied, the conservation properties are recovered as for the unpenalized MLE case.

*2.4. GM component particle sampling*

In physical space, we employ uniform random sampling independently within each cell, which effectively assumes the PDF is constant within each spatial cell. In velocity space (per spatial cell), we employ the ancestral (or forward) sampling technique [29] to generate random samples of a Gaussian mixture. This has the advantage that it allows independent sampling per Gaussian while keeping sampled particle weights identical. To begin, we re-write Eq. 1 as

$$
f(\mathbf{v}) = \sum_{\mathbf{z}} f(\mathbf{z}) f(\mathbf{v}|\mathbf{z}), \tag{16}
$$

where $\mathbf{z}$ is random unit vector of length $K$ (representing the mixture components), with only one non-zero element $z_k = 1$ (chosen randomly) [29]. The identification variable $\mathbf{z}$ has a categorical distribution $f(\mathbf{z})$, and the conditional distribution of $\mathbf{v}$ given $\mathbf{z}$ is a Gaussian. We first draw a sample from $f(\mathbf{z})$, which identifies a Gaussian component $k$ with the probability $\omega_k$. We then draw a sample from the multivariate Gaussian component [30]. In this study, we have used the SPRNG scalable parallel library [31] for random number generation.

Sampling errors in physical space result in violations of Gauss' law (because the accumulated charge density on the mesh will not be identical to the pre-checkpoint state). Similarly, sampling errors in velocity space will break the conservation of mean and variance. Corrections must be made to the sampling procedure to ensure that Gauss' law, momentum and energy are exactly preserved [13]. We discuss these next.

*2.4.1. Preservation of Gauss' law*

In plasmas, Gauss' law is directly related to local charge conservation. After particles are spatially resampled within a cell, the local charge density on the mesh no longer agrees with the pre-checkpoint stage. The local charge density at cell $i$ is given by [32]:

$$\rho_i = \frac{1}{\Delta \mathbf{x}_i} \sum_p q_p \alpha_p S(\mathbf{x}_i - \mathbf{x}_p),$$

with $q_p$ the particle charge, $\alpha_p$ the particle weight, $\mathbf{x}_p$ the particle position within a cell, $\mathbf{x}_i$ the cell center, $\Delta \mathbf{x}_i$ the cell volume, and $S(\mathbf{x})$ a partition-of-unity interpolation kernel (typically a B-spline [32]). Clearly, changes in $\mathbf{x}_i$ will generally result in changes in $\rho_i$, and therefore in Gauss' law, $\nabla \cdot \mathbf{E} = \rho$.

In order to recover the original charge density, we use a technique introduced in Ref. [14] to match the charge density that before checkpointing. The basic idea here is to solve for a slight adjustment of the weight of particles as follows. To be practical, such a weight adjustment is assumed to be uniform within a cell. We begin by assigning a weigh-correction degree of freedom per spatial cell, $\delta A_j$, and define the particle weight correction for all particles in cell $j$ to be equal to $\delta A_j$, i.e.:

$$\delta \alpha_p = \sum_j \delta A_j S_0(\mathbf{x}_j - \mathbf{x}_p), \tag{17}$$

where $S_0(\mathbf{x}_j - \mathbf{x}_p)$ is the zeroth-order B-spline (top-hat) interpolation kernel. The new particle weight is found as:

$$\alpha'_p = \alpha_p + \delta \alpha_p. \tag{18}$$

The weight correction $\delta \alpha_p$ is found by matching the desired charge density $\rho'_i$ (here, using second-order B-splines), i.e.:

$$\rho'_i = \frac{1}{\Delta \mathbf{x}_i} \sum_p q_p \alpha'_p S_2(\mathbf{x}_i - \mathbf{x}_p).$$

Introducing Eqs. 17 and 18 into the last equation, there results:

$$\sum_j \delta A_j \underbrace{\sum_p q_p S_0(\mathbf{x}_j - \mathbf{x}_p) S_2(\mathbf{x}_i - \mathbf{x}_p)}_{M_{ij}} = \Delta \mathbf{x}_i (\rho'_i - \rho_i).$$

The resulting linear system for $\delta A_j$ is a mass-matrix solve, where the matrix is found from contributions from particles to each cell according to stated interpolation rules. By construction, the matrix $\overline{\overline{\mathbf{M}}}$ is sparse, diagonally dominant (because particles in cell $j$ will contribute to that cell the most), and with all positive entries. It is not stiff, and typically the resulting linear system can be converged to round-off in a few iterations.

*2.4.2. Preservation of mean and variance by local projection*

As a result of the mass-matrix Gauss correction step, changes in the particle weight lead to changes in local momentum and energy, breaking strict momentum and energy conservation in each cell. To recover them, we use a well-known projection strategy of particle velocities within each cell proposed by Lemons [13] to correct for momentum and energy errors.

To begin, we follow Ref. [13] and introduce a scaling $\alpha$ and shift $\beta$ for the particle velocity as:

$$\mathbf{v}'_p = \alpha(\mathbf{v}_p + \boldsymbol{\beta}).$$

The parameters $\alpha$ and $\beta$ are determined from the conservation constraints of momentum $\mathbf{p}$ and energy $E$ (which are obtained from the GM):

$$\sum_p \alpha_p \mathbf{v}_p = \sum_p \alpha'_p \mathbf{v}'_p = \mathbf{p} \; ; \; \frac{1}{2} \sum_p \alpha_p v_p^2 = \frac{1}{2} \sum_p \alpha'_p (v'_p)^2 = E.$$

An exact solution $\alpha$ and $\beta$ in terms of $E$ and $\mathbf{p}$ can be found as:

$$\alpha = \sqrt{\frac{2EN'_p - p^2}{2E'N'_p - (p')^2}} \; ; \; \beta = \frac{\mathbf{p} - \alpha \mathbf{p}'}{\alpha N'_p}, \tag{19}$$

where:

$$N'_p = \sum_p \alpha'_p \; ; \; \mathbf{p}' = \sum_p \alpha'_p \mathbf{v}_p \; ; \; E' = \frac{1}{2} \sum_p \alpha'_p v_p^2. \tag{20}$$

Note that even though Schwarz inequality guarantees that $2E'N'_p \geq (p')^2$ (and therefore the denominator in Eq. 19 is always positive definite), the numerator may occasionally become negative (we have seen this when the number of particles is not large enough), and therefore this is a potential failure mode of the approach. When this occurs, there are two options: increase the targeted number of particles for that cell, or forgo the local moment-matching step in that cell.

## 3. Numerical experiments

In this section, we test the proposed CR algorithm using some prototypical test problems, the 1D-1V electrostatic two-stream instability, and the 2D-3V electromagnetic Weibel instability. We perform the simulations with the DPIC code, based on a recently proposed implicit, charge and energy conserving multi-dimensional electromagnetic PIC algorithm [33]. Because of its exact charge- and energy-conserving formulation, DPIC simulations represent a stringent test of the conservation properties (or lack thereof) of the proposed CR algorithm.

*3.1. 1D-1V two-stream electrostatic instability*

The two-stream instability [34] is an electrostatic instability in which two counter-streaming particle beams exchange kinetic and electrostatic energy, and as a result tangle up in to a vortex in phase space [35]. The simulation is performed for $L = 2\pi$ (domain size, in Debye length units), $v_b = \sqrt{3}/2$ (beam

Figure 1: Two-stream instability: Semi-log-scale time history of the electric field energy $E_E$ (top-left), the rms of Gauss' law residual over the whole mesh (top-right), the rms of the residual of the charge conservation equation (bottom-left), and the change of total energy between subsequent time steps (bottom-right). The simulations are obtained without restart, and with GM restart at $t = 10$ (in normalized units) with and without Lemons moment matching.

speed, in electron thermal speed units), $N_x = 32$ (number of cells), $N_p = 156$ (number of particles per cell), $\Delta t = 0.2$ (time step in inverse plasma frequency units), with periodic boundary conditions. Figure 1 shows the root-mean-square (rms) of the charge conservation equation residual $(\partial_t \rho + \nabla \cdot \mathbf{j})$ over the mesh, the electric-field energy $E_E = \sum_i \frac{E_i^2}{2}$, the total energy error between subsequent timesteps, $|\mathcal{E}^{n+1} - \mathcal{E}^n|$, with $\mathcal{E}$ the total sum of particle and electric-field energy, and the rms of the residual of Gauss' law, $\nabla \cdot E - \rho$. The plot compares the unrestarted run with two GM-restarted ones at $t = 10$ (mid/late linear stage), with and without Lemons' moment matching. The results show exact conservation of charge for all cases, also for energy except for the case where Lemons matching was not used (which results in a large energy conservation error right after restart), and excellent preservation of Gauss' law (commensurate with the nonlinear tolerance). They also show excellent agreement in the temporal evolution of the electrostatic field energy for all cases. For this run, the GM algorithm is started with 8 Gaussian's per cell, resulting in an average number of Gaussians per cell of 2, and therefore to an average compression ratio of about 75.

A comparison between 1D-1V phase-space plots between unrestarted (left) and GM-restarted (right) runs is shown in Fig. 2. It can be appreciated that the GM-restarted phase-space plot (right) captures all phase-space features present in the unrestarted case (left), except for a bit of beam-spread in the particles (which is generated by the random uniform spatial initialization per cell in the GM-restarted case).

Figure 2: Two-stream instability: Phase-space comparison at three different times (green: $t = 0$, blue: $t = 14.0$, red: $t = 19.4$) between the unrestarted case (left) and the GM-restarted one (right).

### 3.2. 2D-3V Weibel electromagnetic instability

Next we test with Weibel instability, which is a electromagnetic instability in a plasma with anisotropic temperatures [36]. Unless otherwise specified, simulations below are performed in a 2D domain $10d_e \times 10d_e$ (where $d_e$ is the electron skin depth), with $16 \times 16$ cells, $\Delta t = 1$ (in inverse plasma frequency units), with doubly periodic boundary conditions. A temperature anisotropy is set up for both electrons and ions with $v_{thx} = 0.1$, and $v_{thy,z} = 0.3$ in speed-of-light units. The mass ratio is set to be $m_i/m_e = 1836$. We initialize the simulation with a $\delta$-function perturbation in the particle velocities of $10^{-3}$, as described in [33]. Note that the code employed has assumed Darwin approximation [33], which is non-relativistic, and does not admit any light wave propagation in the system.

Figure 3 shows similar time histories as in Fig. 1 but with the magnetic energy $E_B = \sum_i B_i^2/2$, obtained with $N_p = 128$ particles per cell, and with and without restart at $t = 20$ in normalized time units (late in the linear phase). It is apparent that conservation properties are preserved before and after restart (except for energy without Lemons projection, as expected), and that the GM restart quality is quite good, even with this relatively small number of particles per cell. Increasing the number of particles per cell improves the agreement, as it is shown in Fig. 4. The initial number of Gaussians per cell is 8, leading to an average number of Gaussians per cell of 1.8, 2.1, 3.3 for 128, 512, and 1024 particles per cell, respectively, implying a compression ratio of 70, 240, and 310.

The ability of the approach to deal with particles with different weights is shown in Fig. 5, which depicts similar time histories as before but now for the unrestarted and twice-GM-restarted (at $t = 15$ and $t = 30$) 2D Weibel instability with 512 particles per cell. The simulation begins with identical particles, but their weights develop differences due to the density mass-matrix solve at the first GM restart. The second GM restart is therefore performed with non-identical particles. Agreement between restarted and unrestarted time histories is very good throughout the simulation, and demonstrates the ability of the method to deal with particles with arbitrary weight.

Figure 3: 2D Weibel instability with $N_p = 128$: Semi-log-scale time history of the magnetic field energy $E_B$ (top-left), the rms of the Gauss' law residual over the whole mesh (top-right), the rms of residual of the charge conservation equation (bottom-left), and the change of total energy between subsequent time steps (bottom-right). The simulations are obtained without restart, and with GM restart at $t = 20$ (in normalized units) with and without Lemons moment matching.



Figure 4: 2D Weibel instability: Semi-log-scale time history of the magnetic field energy with $N_p = 128$ (left), 512 (center), and 1024 (right). The simulations are obtained without restart, and with GM restart at $t = 20$ (in normalized units) with and without Lemons moment matching.

Figure 5: Time histories of the same quantities as in Fig. 3, comparing the twice-restarted 2D Weibel instability with 512 particles per cell vs. the unrestarted result, demonstrating the ability of the method to deal with particles of arbitrary weight.

### 3.3. Particle remapping using EM-GM for noise reduction (variance control)

The central goal of machine-learning algorithms is not only to provide a goodness-of-fit to the data, but also to be able to generalize. The implication in the context of GM is that the estimated continuum PDF may be able distinguishing between noise and signal, and, if so, provide a measure of noise reduction (i.e., variance control), such that the GM PDF, once resampled, may lead to an improved PIC solution vs. the unrestarted one. The subject of noise control in PIC algorithms has received significant attention recently [37, 38, 39], but it has mostly been circumscribed to the remapping of the particle PDF via interpolation to a (semi-)structured phase-space mesh (i.e., bins), and subsequent resampling within bins. Some of these approaches [39] explicitly embed arbitrary moment conservation in their formulation, which is a desirable property. However, to our knowledge, the use of Gaussian-mixture techniques for this purpose remains unexplored.

Here, we provide anecdotal evidence that particle remapping using the GM PDF reconstruction proposed here actually leads to an improvement in PIC solution quality, suggesting that a thorough exploration of this subject is worthwhile (and will be the subject of future work). For our demonstration, we choose a Weibel instability in a 1D domain of size $\pi$ (in $d_e$ units), with $\Delta t = 1$ (in inverse plasma frequency units), and periodic boundary conditions. The temperature anisotropy and mass ratio is the same as in the previous Weibel example. We initialize the simulation with a $\delta$-function perturbation in the particle velocities of $10^{-2}$.

13

Figure 6: 1D Weibel instability comparison with and without restart at $t = 15$ for low resolution (left), high resolution (center), and the comparison of the two (right). Note that, for the right plot, a time shift ($t \leftarrow t + 12$) has been applied to the low-resolution histories to facilitate a meaningful assessment.



Figure 7: Performance of EM algorithm (relative change of log-likelihood vs iteration) with and without AA and K-means initialization at cell (1,3) of 2D Weibel problem at $t = 20$ with different particle resolutions.

Figure 6 shows a comparison of the magnetic-field energy evolution between unrestarted and GM-restarted PIC simulations, for a low-resolution case ($N_p = 1000$, $N_x = 32$, left) and a high-resolution one ($N_p = 4000$, $N_x = 128$, center). Both are restarted at $t = 15$. For the low resolution case (Fig. 6-left), one can appreciate a relatively big difference in the evolution of the magnetic-field energy between the unrestarted and GM-restarted simulations, especially when it enters the nonlinear stage. The magnetic field energy is higher in the GM-restarted simulation, and there are also some phase differences in the nonlinear oscillation. As we reduce the grid size and increase the number of particles per cell, however, the history of the magnetic-field energy agree much better between unrestarted and GM-restarted solutions (Fig. 6-center), indicating that the simulation is converging. More interestingly, when one compares the low-resolution simulations with the high-resolution ones (Fig. 6-right), it is apparent that the GM-restarted low-resolution solution is much closer to the high-resolution result than the unrestarted low-resolution one. It follows that, for a given resolution, the GM-restarted simulation is able to achieve a more accurate B-field nonlinear saturation energy level than the unrestarted one, suggesting that the generalization capability rooted in the unsupervised machine-learning algorithm is in fact at play.

### 3.4. On the acceleration of convergence of the EM-GM nonlinear algorithm

As discussed previously, the EM-GM algorithm is guaranteed to converge [10], but performance can be slow. There has been recent work trying to accelerate the convergence of the EM algorithm, both by improv-

14

ing the initialization of the iteration (e.g., using K-means [40]), or by improving the Picard iteration itself (e.g., by using Anderson Acceleration (AA) [41, 40] or by advanced conjugate search direction algorithms [42, 43]). In this study, we have implemented the K-means initialization and the Anderson Acceleration algorithm in the *standard* (non-adaptive) EM-GM algorithm, and tested its impact using Weibel instability data.

A word is in order about our AA implementation for EM-GM, which to our knowledge is new. We have included in the residual all degrees of freedom for all Gaussians, namely, all weights $\omega_k$, means $\boldsymbol{\mu}_k$, and second-moment matrices (i.e., $\overline{\overline{\mathbf{M}}}_{2,k} = \int d\mathbf{v}\,\mathbf{v}\mathbf{v}^T f_k$, instead of covariance matrices, $\overline{\overline{\boldsymbol{\Sigma}}}_k = \int d\mathbf{v}(\mathbf{v} - \boldsymbol{\mu}_k)(\mathbf{v} - \boldsymbol{\mu}_k)^T f_k = \overline{\overline{\mathbf{M}}}_{2,k} - \boldsymbol{\mu}_k\boldsymbol{\mu}_k^T$). The latter choice is motivated by the fact that $\overline{\overline{\mathbf{M}}}_{2,k}$ is linear in the mixture components (i.e., the second moment of a linear combination of Gaussians is the linear combination of the second moments of each individual Gaussian), whereas $\overline{\overline{\boldsymbol{\Sigma}}}_k$ is nonlinear, and therefore the former are better suited for acceleration based on a linear combination of past residuals. It also has potential advantages for the preservation of positivity of $\overline{\overline{\boldsymbol{\Sigma}}}_k$ (see below). It is important to note that, unlike EM, the standard AA algorithm does not conserve moments, and does not guarantee that covariance matrices remain positive definite. To fix conservation, we apply a *standard* EM step after the AA iteration (as was done for the penalized EM algorithm for the same reason). To fix positivity, we currently reset the AA iteration after an indefinite covariance matrix is detected and revert back to a standard EM step. An alternate approach (enabled by our choice to accelerate $\overline{\overline{\mathbf{M}}}_{2,k}$ instead of $\overline{\overline{\boldsymbol{\Sigma}}}_k$) would be to guarantee that the Anderson mixing coefficients remain positive (i.e., that the linear combination of residuals in AA remain convex). Recent studies promote this as a viable globalization procedure for AA [44], and this strategy will be explored in future work.

In Fig. 7, we report on the impact of these strategies on the convergence rate of the for cell (3,1) of the 2D Weibel test using 128, 512, and 1024 particles per cell. The plots demonstrate that AA and K-means initialization (both independently and combined) result in a significant speedup of the rate of convergence of the non-adaptive EM-GM algorithm, which is more noticeable with increasing number of particles per cell. It is apparent, however, that AA is much more effective in accelerating the EM-GM convergence than K-means. From these plots, the speedup between the unaccelerated random initialization case to the AA accelerated K-means initialization case is of more than an order of magnitude. These speedups are in fact representative of performance in most cells for the Weibel test problem. We are currently exploring ways of generalizing these strategies for the *adaptive* EM-GM algorithm, and these will be reported in a future publication.

## 4. Discussion and summary

We have proposed a checkpoint-restart strategy for PIC algorithms based on unsupervised machine-learning strategies using Gaussian Mixture models. The Gaussian components are found adaptively using a penalized Maximum Likelihood Estimate, solved by an Expectation-Maximization procedure. For the

numerical tests presented, the approach has demonstrated significant compression potential (of several orders of magnitude) without loss of physical fidelity (as demonstrated by actual restarted PIC simulations). The latter is facilitated by the exact preservation of charge, momentum, and energy in both compression and reconstruction stages, and by the fact that GMM provides an optimal continuum reconstruction of the PDF represented by the particles. Key to the fidelity of the approach (particularly if many CRs are performed) is the use of a mass-matrix procedure to match the density profile on the mesh exactly, and a projection step to enforce conservation properties after particle resampling. Our numerical experiments not only demonstrate that the approach successfully restarts both electrostatic and electromagnetic PIC simulations (with strict conservation of both charge and energy exactly, and which therefore represent a stringent test of the method), but also suggest that a periodic GM particle remap may in fact improve the quality of PIC solutions. This point, which is anecdotal in this study, suggests the possibility of machine-learning variance reduction in particle methods, and will be investigated further in future work. Finally, we have proposed a simple implementation strategy for Anderson Acceleration in the *non-adaptive* EM-GM algorithm that results in convergence speedups of more than an order of magnitude, while strictly preserving the positivity of the covariance matrices of the mixture. Beyond CR and particle remapping for enhanced solution quality, we note that the approach outlined in this study enables straightforwardly particle redistribution over the computational domain, to facilitate various performance goals such as load balancing and particle-number control. This will also be the subject of future studies.

**Appendix A. Derivation of Penalized Likelihood Function**

We begin by making use of the exact decomposition [45]:

$$\ln p(\mathbf{X}|K) = L(q) + KL(q||p), \tag{A.1}$$

where

$$L(q) = \int q(\boldsymbol{\theta}) \ln \frac{p(\mathbf{X}|\boldsymbol{\theta}, K) p(\boldsymbol{\theta}|K)}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}, \tag{A.2}$$

$$KL(q||p) = \int q(\boldsymbol{\theta}) \ln \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{X}, K)} d\boldsymbol{\theta}, \tag{A.3}$$

Note that the decomposition holds for an arbitrary positive-definite distribution $q(\boldsymbol{\theta})$. Since the Kullback-Leibler divergence $KL(q||p)$ is always greater or equal to zero [19], $L(q)$ is a lower bound of the log-marginal likelihood. In fact, maximizing $L(q)$ is equivalent to maximizing $\ln p(\mathbf{X}|K)$ [45]. Various forms of $q(\boldsymbol{\theta})$ can be adopted. For instance, variational Bayesian methods assume that $q(\boldsymbol{\theta})$ factorizes over subsets $\{\boldsymbol{\theta}_i\}$, i.e., $q(\boldsymbol{\theta}) = \Pi_i q_i(\boldsymbol{\theta}_i)$ [46]. Here, $q(\boldsymbol{\theta})$ is assumed to be a uniform distribution in a small interval $(\mathbf{a}, \mathbf{b})$ around a point of $\boldsymbol{\theta}$, i.e., $q(\boldsymbol{\theta}) = 1/\Delta$, with $\Delta \equiv \int_{\mathbf{a}}^{\mathbf{b}} d\boldsymbol{\theta}$ the volume of a $d$-dimensional hypercube in the space of parameter $\boldsymbol{\theta}$. Equation A.2 then becomes

$$L(\Delta) = \frac{1}{\Delta} \int_{\mathbf{a}}^{\mathbf{b}} \ln\left[p(\boldsymbol{\theta}|K)\Delta\right] d\boldsymbol{\theta} + \frac{1}{\Delta} \int_{\mathbf{a}}^{\mathbf{b}} \ln\left[p(\mathbf{X}|\boldsymbol{\theta}, K)\right] d\boldsymbol{\theta}. \tag{A.4}$$

In the context of information theory, the maximum of Eq. A.4 is equivalent to the shortest message length that the data can communicate [19]. The idea is that the model with the minimum message length (thus so-called MML) should be preferred. The message length, defined as $l(x) = -\ln P(x)$, where $P(x)$ is the probability of an event $x$, is a measure of the information content of the event $x$ [19]. For a continuous PDF $p(x)$, $l(x) = -\ln[p(x)dx]$, with $dx$ a small interval around $x$. It is clear from this perspective that the first term on the right-hand-side (rhs) of Eq. A.4 corresponds to the message length of $\boldsymbol{\theta}$, and $\Delta$ denotes a discretization of $\boldsymbol{\theta}$ (which may be thought as the finite precision of $\boldsymbol{\theta}$). The finite precision of $\boldsymbol{\theta}$ has a major effect on "communicating" the message length of the data. Expectation is taken with respect to the assumed uniform distribution $q(\boldsymbol{\theta})$ over a small interval $(\mathbf{a}, \mathbf{b})$, and an optimum $\Delta$ can be found by maximizing Eq. A.4.

We next rewrite the log-marginal likelihood function (Eq. A.4) as:

$$L(\Delta_\xi) = \frac{1}{\Delta_\xi} \int_{\boldsymbol{\alpha}}^{\boldsymbol{\beta}} \ln\left[p(\boldsymbol{\xi}|K, \boldsymbol{\omega})p(\boldsymbol{\omega})\Delta_\xi\right] d\boldsymbol{\xi} + \frac{1}{\Delta_\xi} \int_{\boldsymbol{\alpha}}^{\boldsymbol{\beta}} \ln\left[p(\mathbf{X}|\boldsymbol{\xi}, K, \boldsymbol{\omega})\right] d\boldsymbol{\xi}. \tag{A.5}$$

where we have made a variable transformation $\boldsymbol{\theta} = \Lambda^{-1/2}U^T\boldsymbol{\xi}$ with Jacobian $J = |\partial\boldsymbol{\theta}/\partial\boldsymbol{\xi}|$. Here $U$ is a $d \times d$ orthogonal matrix with columns given by eigenvectors and $\Lambda$ is a $d \times d$ diagonal matrix with elements of eigenvalues of the observed Fisher information matrix, $I_p = -\left.\frac{\partial^2 \ln p(\mathbf{X}|\boldsymbol{\theta}, K, \boldsymbol{\omega})}{\partial\theta^2}\right|_{\hat{\boldsymbol{\theta}}}$. Here we have assumed that the Hessian matrix is a negative semidefinite (e.g., when $\ln p$ is concave), so that we can write $I_p = U\Lambda U^T$, where $U$ is an orthogonal matrix. It follows that $J = |I_p|^{-\frac{1}{2}}$. Using the chain rule [47] we find that $I_p(\boldsymbol{\xi}) = (\partial\boldsymbol{\theta}/\partial\boldsymbol{\xi})^T I_p (\partial\boldsymbol{\theta}/\partial\boldsymbol{\xi}) = \mathbb{1}$. A truncated Taylor expansion with respect to the center (denoted as $\tilde{\boldsymbol{\xi}}$) of $\Delta_\xi$ is typically employed to approximate the log-likelihood $\ln p(\mathbf{X}|\boldsymbol{\xi}, K, \boldsymbol{\omega})$:

$$\ln p(\mathbf{X}|\boldsymbol{\xi}, K, \boldsymbol{\omega}) \simeq \ln p(\mathbf{X}|\tilde{\boldsymbol{\xi}}, K, \boldsymbol{\omega}) + (\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}}) \cdot \left.\frac{\partial\ln p}{\partial\boldsymbol{\xi}}\right|_{\tilde{\boldsymbol{\xi}}} + \frac{1}{2}(\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}})^T \cdot (\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}}). \tag{A.6}$$

Substituting Eq. A.6 into Eq. A.5 results

$$L(\tilde{\boldsymbol{\xi}}, \Delta_\xi) = \ln\left[p(\mathbf{X}|\tilde{\boldsymbol{\xi}}, K, \boldsymbol{\omega})p(\tilde{\boldsymbol{\xi}}|K, \boldsymbol{\omega})p(\boldsymbol{\omega})\Delta_\xi\right] - \frac{d}{24}\Delta_\xi^{2/d}, \tag{A.7}$$

where we have used $\frac{1}{\Delta_\xi}\int_{\boldsymbol{\alpha}}^{\boldsymbol{\beta}}(\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}})d\boldsymbol{\xi} = 0$, and $\frac{1}{\Delta_\xi}\int_{\boldsymbol{\alpha}}^{\boldsymbol{\beta}}(\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}})^T \cdot (\boldsymbol{\xi} - \tilde{\boldsymbol{\xi}})d\boldsymbol{\xi} = \frac{d}{12}\Delta_\xi^{2/d}$, both integrated over the volume $\Delta_\xi$ (a $d$-dimensional hypercube), and assuming that all the other terms are constant within $\Delta_\xi$. By

setting $\frac{\partial L}{\partial \Delta_\xi} = 0$, $\Delta_\xi = (12)^{d/2}$ is found to maximize Eq. A.7. Substituting $\Delta_\xi = (12)^{d/2}$ into Eq. A.7 yields:

$$L(\tilde{\boldsymbol{\theta}}) = \ln\left[p(\mathbf{X}|\tilde{\boldsymbol{\theta}}, K, \omega)p(\tilde{\boldsymbol{\theta}}|K, \omega)p(\omega)\right] - \frac{1}{2}\ln|I_p| - \frac{d}{2}(1 - \ln 12), \tag{A.8}$$

where we have also used $p(\mathbf{X}|\tilde{\boldsymbol{\theta}}, K)p(\tilde{\boldsymbol{\theta}}|K)J = p(\mathbf{X}|\tilde{\boldsymbol{\xi}}, K)p(\tilde{\boldsymbol{\xi}}|K)$ due to the variable transformation [48]. The negative of Equation A.8 is the so-called MML criterion [49]:

$$\textit{Message Length} = -\ln\left[p(\mathbf{X}|\tilde{\boldsymbol{\theta}}, K, \omega)p(\tilde{\boldsymbol{\theta}}|K, \omega)p(\omega)\right] + \frac{1}{2}\ln|I_p| + \frac{d}{2}(1 - \ln 12). \tag{A.9}$$

Note that $\Delta_\xi$ related terms group into the last term, which is in general not that important (see below).

We must further simplify Eq. A.8 because of the difficulties in selecting prior distributions [50] and calculating the Fisher information. We start with noting that $I_p = -\left.\frac{\partial^2 \sum_{i=1}^{N} \ln p(\mathbf{x}_i|\boldsymbol{\theta}, K)}{\partial \theta^2}\right|_{\tilde{\boldsymbol{\theta}}} = -\sum_{i=1}^{N} \left.\frac{\partial^2 \ln p(\mathbf{x}_i|\boldsymbol{\theta}, K)}{\partial \theta^2}\right|_{\tilde{\boldsymbol{\theta}}} \simeq N\mathcal{I}$ where $\mathcal{I} = -E_{\boldsymbol{\theta}}(\partial^2 \ln p(\mathbf{x}|\boldsymbol{\theta}, K)/\partial \theta^2)$ is the Fisher information matrix (FIM), and the expectation is taken with the mixture PDF [8]. To proceed, only an upper bound of $|\mathcal{I}|$, the complete FIM $\mathcal{I}_c$ [51], is considered [52],

$$\mathcal{I}_c = \text{Blockdiag}(\omega_1 I_1, ..., \omega_K I_K, I_\omega), \tag{A.10}$$

where $I_k$ is a $T \times T$ FIM, with $T = \frac{1}{2}D(D + 3)$ and recall that $D$ is the dimension of $\boldsymbol{\mu}$. Note that $d \times d$ is the dimension of $I$, where $d = KT + K - 1$ is the total number of parameters. The minus one is due to the constraint that $\sum_i \omega_i = 1$. The second term on the rhs of Eq. A.8 may be written as $-\frac{1}{2}\ln(N^d|\mathcal{I}_c|) = -\frac{1}{2}\ln[N^d \prod_{i=1}^{K}(\omega_i^T|I_i|)|I_\omega|]$. With the above approximations, Eq. A.8 yields:

$$L(\tilde{\boldsymbol{\theta}}) = \ln\left(\frac{p(\tilde{\boldsymbol{\theta}}|K)}{\sqrt{|I_\omega|}\prod_{i=1}^{K}\sqrt{|I_i|}}\right) + \ln p(\mathbf{X}|\tilde{\boldsymbol{\theta}}, K) - \frac{d}{2}\ln N - \frac{T}{2}\sum_{i=1}^{K}\ln \omega_i - \frac{d}{2}(1 - \ln 12). \tag{A.11}$$

If we choose independent priors, i.e., $p(\tilde{\boldsymbol{\theta}}|K) = p(\boldsymbol{\omega})\prod_{i=1}^{K} p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, and adopt Jeffreys' prior for $\boldsymbol{\omega}$ and for each $(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ [53, 9], we obtain

$$L(\boldsymbol{\omega}, \boldsymbol{\mu}, \overline{\overline{\boldsymbol{\Sigma}}}) = \ln\left(p(\mathbf{X}|\boldsymbol{\omega}, \boldsymbol{\mu}, \overline{\overline{\boldsymbol{\Sigma}}}, K)\right) - \frac{d}{2}\ln N - \frac{T}{2}\sum_{i=1}^{K}\ln(\omega_i), \tag{A.12}$$

after dropping some constants and $\sim O(d)$ terms (more specifically the condition for dropping the last term of Eq. A.8 is $N \gg 1$ which is typically the case). We have arrived at a simple penalized likelihood function, Eq. A.12. It is worth noting that the MML estimator of maximizing the likelihood Eq. A.12 is invariant under variable transformation, or re-parameterization of $(\boldsymbol{\mu}, \overline{\overline{\boldsymbol{\Sigma}}})$ and $\boldsymbol{\omega}$. This is due to the invariance property of maximum likelihood estimators to arbitrary transformations of the parameters of likelihood function [48].

## References

[1] J. A. Kahle, J. Moreno, and D. Dreps, "2.1 summit and sierra: Designing ai/hpc supercomputers," in *2019 IEEE International Solid-State Circuits Conference-(ISSCC)*, pp. 42–43, IEEE, 2019.

[2] E. B. Nightingale, J. R. Douceur, and V. Orgovan, "Cycles, cells and platters: an empirical analysisof hardware failures on a million consumer pcs," in *Proceedings of the sixth conference on Computer systems*, pp. 343–356, 2011.

[3] R.-T. Liu and Z.-N. Chen, "A large-scale study of failures on petascale supercomputers," *Journal of computer science and technology*, vol. 33, no. 1, pp. 24–41, 2018.

[4] E. Rojas, E. Meneses, T. Jones, and D. Maxwell, "Analyzing a five-year failure record of a leadership-class supercomputer," in *2019 31st International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*, pp. 196–203, IEEE, 2019.

[5] D. Dauwe, S. Pasricha, A. A. Maciejewski, and H. J. Siegel, "An analysis of resilience techniques for exascale computing platforms," in *2017 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pp. 914–923, IEEE, 2017.

[6] Z. Miao, J. Calhoun, and R. Ge, "Energy analysis and optimization for resilient scalable linear systems," in *2018 IEEE International Conference on Cluster Computing (CLUSTER)*, pp. 24–34, IEEE, 2018.

[7] J. F. Lofstead, S. Klasky, K. Schwan, N. Podhorszki, and C. Jin, "Flexible io and integration for scientific codes through the adaptable io system (adios)," in *Proceedings of the 6th international workshop on Challenges of large applications in distributed environments*, pp. 15–24, 2008.

[8] G. J. McLachlan and D. Peel, *Finite mixture models*. John Wiley & Sons, 2004.

[9] M. A. Figueiredo and A. K. Jain, "Unsupervised selection and estimation of finite mixture models," in *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, vol. 2, pp. 87–90, IEEE, 2000.

[10] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the royal statistical society. Series B (methodological)*, pp. 1–38, 1977.

[11] C. S. Wallace, *Statistical and inductive inference by minimum message length*. Springer Science & Business Media, 2005.

[12] J. Behboodian, "On a mixture of normal distributions," *Biometrika*, vol. 34, no. 57 Part 1, pp. 215–217, 1970.

[13] D. S. Lemons, D. Winske, W. Daughton, and B. Albright, "Small-angle coulomb collision model for particle-in-cell simulations," *Journal of Computational Physics*, vol. 228, no. 5, pp. 1391–1403, 2009.

[14] D. Burgess, D. Sulsky, and J. Brackbill, "Mass matrix formulation of the flip particle-in-cell method," *Journal of Computational Physics*, vol. 103, no. 1, pp. 1–15, 1992.

[15] R. Dupuis, M. V. Goldman, D. L. Newman, J. Amaya, and G. Lapenta, "Characterizing magnetic reconnection regions using gaussian mixture models on particle velocity distributions," *The Astrophysical Journal*, vol. 889, no. 1, p. 22, 2020.

[16] K. J. Bowers, B. G. Devolder, L. Yin, and T. J. Kwan, "A maximum likelihood method for linking particle-in-cell and monte-carlo transport simulations," *Computer physics communications*, vol. 164, no. 1-3, pp. 311–317, 2004.

[17] D. J. Larson and C. V. Young, "A finite mass based method for vlasov–poisson simulations," *Journal of Computational Physics*, vol. 284, pp. 171–185, 2015.

[18] B. S. Everitt, "Finite mixture distributions," *Wiley StatsRef: Statistics Reference Online*, 2014.

[19] D. J. MacKay and D. J. Mac Kay, *Information theory, inference and learning algorithms*. Cambridge university press, 2003.

[20] I. J. Good, *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. MIT Press, 1965.

[21] J. Rousseau and K. Mengersen, "Asymptotic behaviour of the posterior distribution in overfitted mixture models," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 5, pp. 689–710, 2011.

[22] Z. Zivkovic, "Improved adaptive gaussian mixture model for background subtraction," in *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 2, pp. 28–31, IEEE, 2004.

[23] K. Tu, "Modified dirichlet distribution: Allowing negative parameters to induce stronger sparsity," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 1986–1991, 2016.

[24] V. Hasselblad, "Estimation of parameters for a mixture of normal distributions," *Technometrics*, vol. 8, no. 3, pp. 431–444, 1966.

[25] J.-L. Gauvain and C.-H. Lee, "Maximum a posteriori estimation for multivariate gaussian mixture observations of markov chains," *IEEE transactions on speech and audio processing*, vol. 2, no. 2, pp. 291–298, 1994.

[26] R. A. Redner and H. F. Walker, "Mixture densities, maximum likelihood and the em algorithm," *SIAM review*, vol. 26, no. 2, pp. 195–239, 1984.

[27] G. Celeux, S. Chrétien, F. Forbes, and A. Mkhadri, "A component-wise em algorithm for mixtures," *Journal of Computational and Graphical Statistics*, vol. 10, no. 4, pp. 697–712, 2001.

[28] J. K. Blitzstein and J. Hwang, *Introduction to probability*. Crc Press, 2019.

[29] C. M. Bishop, *Pattern recognition and machine learning*. springer, 2006.

[30] Y. L. Tong, *The multivariate normal distribution*. Springer Science & Business Media, 2012.

[31] M. Mascagni and A. Srinivasan, "Algorithm 806: Sprng: A scalable library for pseudorandom number generation," *ACM Transactions on Mathematical Software (TOMS)*, vol. 26, no. 3, pp. 436–461, 2000.

[32] C. K. Birdsall and A. B. Langdon, *Plasma physics via computer simulation*. CRC press, 2004.

[33] G. Chen and L. Chacon, "A multi-dimensional, energy-and charge-conserving, nonlinearly implicit, electromagnetic vlasov–darwin particle-in-cell algorithm," *Computer Physics Communications*, vol. 197, pp. 73–87, 2015.

[34] M. A. Lampert, "Plasma oscillations at extremely high frequencies," *Journal of Applied Physics*, vol. 27, no. 1, pp. 5–11, 1956.

[35] K. Roberts and H. L. Berk, "Nonlinear evolution of a two-stream instability," *Physical Review Letters*, vol. 19, no. 6, p. 297, 1967.

[36] E. S. Weibel, "Spontaneously growing transverse waves in a plasma due to an anisotropic velocity distribution," *Physical Review Letters*, vol. 2, no. 3, p. 83, 1959.

[37] B. Wang, G. H. Miller, and P. Colella, "A particle-in-cell method with adaptive phase-space remapping for kinetic plasmas," *SIAM Journal on Scientific Computing*, vol. 33, no. 6, pp. 3509–3537, 2011.

[38] A. Myers, P. Colella, and B. V. Straalen, "A 4th-order particle-in-cell method with phase-space remapping for the vlasov–poisson equation," *SIAM Journal on Scientific Computing*, vol. 39, no. 3, pp. B467–B485, 2017.

[39] D. Faghihi, V. Carey, C. Michoski, R. Hager, S. Janhunen, C.-S. Chang, and R. Moser, "Moment preserving constrained resampling with applications to particle-in-cell methods," *Journal of Computational Physics*, vol. 409, p. 109317, 2020.

[40] J. H. Plasse, "The em algorithm in multivariate gaussian mixture models using anderson acceleration," 2013.

[41] H. F. Walker and P. Ni, "Anderson acceleration for fixed-point iterations," *SIAM Journal on Numerical Analysis*, vol. 49, no. 4, pp. 1715–1735, 2011.

[42] Y. He and C. Liu, "The dynamic "expectation–conditional maximization either" algorithm," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 74, no. 2, pp. 313–336, 2012.

[43] W. Xiang, A. Karfoul, C. Yang, H. Shu, and R. L. B. Jeannès, "An exact line search scheme to accelerate the em algorithm: Application to gaussian mixture models identification," *Journal of Computational Science*, p. 101073, 2020.

[44] X. Chen and C. Kelley, "Convergence of the ediis algorithm for nonlinear equations," *SIAM Journal on Scientific Computing*, vol. 41, no. 1, pp. A365–A379, 2019.

[45] R. M. Neal and G. E. Hinton, "A view of the em algorithm that justifies incremental, sparse, and other variants," in *Learning in graphical models*, pp. 355–368, Springer, 1998.

[46] A. Corduneanu and C. M. Bishop, "Variational bayesian model selection for mixture distributions," in *Artificial intelligence and Statistics*, vol. 2001, pp. 27–34, Morgan Kaufmann Waltham, MA, 2001.

[47] M. J. Schervish, *Theory of statistics*. Springer Science & Business Media, 2012.

[48] G. Casella and R. L. Berger, *Statistical inference*, vol. 2. Duxbury Pacific Grove, CA, 2002.

[49] A. D. Lanterman, "Schwarz, wallace, and rissanen: Intertwining themes in theories of model selection," *International statistical review*, vol. 69, no. 2, pp. 185–212, 2001.

[50] R. E. Kass and L. Wasserman, "The selection of prior distributions by formal rules," *Journal of the American Statistical Association*, vol. 91, no. 435, pp. 1343–1370, 1996.

[51] D. M. Titterington, A. F. Smith, and U. E. Makov, *Statistical analysis of finite mixture distributions*. Wiley,, 1985.

[52] A. M. Raim, N. K. Neerchal, and J. G. Morel, "An approximation to the information matrix of exponential family finite mixtures," *Annals of the Institute of Statistical Mathematics*, vol. 69, no. 2, pp. 333–364, 2017.

[53] J. Bernardo and F. Girón, "A bayesian analysis of simple mixture problems," *Bayesian statistics*, vol. 3, no. 3, pp. 67–78, 1988.