

## UvA-DARE (Digital Academic Repository)

### Multiple linear regression and thermodynamic fluctuations are equivalent for computing thermodynamic derivatives from molecular simulation

Rahbari, A.; Josephson, T.R.; Sun, Y.; Moulton, O.A.; Dubbeldam, D.; Siepmann, J.I.; Vlugt, T.J.H.

**DOI**

[10.1016/j.fluid.2020.112785](https://doi.org/10.1016/j.fluid.2020.112785)

**Publication date**

2020

**Document Version**

Final published version

**Published in**

Fluid Phase Equilibria

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

Rahbari, A., Josephson, T. R., Sun, Y., Moulton, O. A., Dubbeldam, D., Siepmann, J. I., & Vlugt, T. J. H. (2020). Multiple linear regression and thermodynamic fluctuations are equivalent for computing thermodynamic derivatives from molecular simulation. *Fluid Phase Equilibria*, 523, [112785]. <https://doi.org/10.1016/j.fluid.2020.112785>

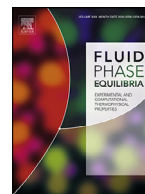
**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

*UvA-DARE is a service provided by the library of the University of Amsterdam (<https://dare.uva.nl>)*



# Multiple linear regression and thermodynamic fluctuations are equivalent for computing thermodynamic derivatives from molecular simulation

Ahmadreza Rahbari<sup>a</sup>, Tyler R. Josephson<sup>b</sup>, Yangzesheng Sun<sup>b</sup>, Othonas A. Moulton<sup>a</sup>, David Dubbeldam<sup>c</sup>, J. Ilja Siepmann<sup>b</sup>, Thijs J.H. Vlugt<sup>a,\*</sup>

<sup>a</sup> Engineering Thermodynamics, Process & Energy Department, Faculty of Mechanical, Maritime and Materials Engineering, Delft University of Technology, Leeghwaterstraat 39, Delft, 2628CB, the Netherlands

<sup>b</sup> Department of Chemistry and Chemical Theory Center, University of Minnesota, 207 Pleasant Street SE, Minneapolis, MN 55455-0431, USA

<sup>c</sup> Van't Hoff Institute for Molecular Sciences, University of Amsterdam, Science Park 904, 1098XH Amsterdam, the Netherlands

## ARTICLE INFO

### Article history:

Received 16 March 2020

Revised 6 July 2020

Accepted 9 August 2020

Available online 11 August 2020

### Keywords:

Molecular simulation

Partial molar properties

Open ensembles thermodynamic fluctuations

Linear regression

## ABSTRACT

Partial molar properties are of fundamental importance for understanding properties of non-ideal mixtures. Josephson and co-workers (Mol. Phys. 2019, 117, 3589–3602) used least squares multiple linear regression to obtain partial molar properties in open constant-pressure ensembles. Assuming composition-independent partial molar properties for the narrow composition range encountered throughout simulation trajectories, we rigorously prove the equivalence of two approaches for computing thermodynamic derivatives in open ensembles of an  $n$ -component system: (1) multiple linear regression, and (2) thermodynamic fluctuations. Multiple linear regression provides a conceptually simple and computationally efficient way of computing thermodynamic derivatives for multicomponent systems. We show that in the reaction ensemble, the reaction enthalpy can be computed directly by simple multiple linear regression of the enthalpy as a function of the number of reactant molecules. Non-linear regression and a Gaussian process model taking into account the compositional dependence of partial molar properties further support that multiple linear regression captures the correct physics.

© 2020 The Author(s). Published by Elsevier B.V.

This is an open access article under the CC BY license. (<http://creativecommons.org/licenses/by/4.0/>)

## 1. Introduction

Partial molar properties are important quantities for describing multicomponent non-ideal mixtures [1–3]. The partial molar property of component  $i$  in a mixture is defined as  $\bar{x}_i = (\partial X / \partial n_i)_{T,P,n_{j \neq i}}$ , in which  $X$  is the corresponding extensive property of the mixture and  $n_i$  is the number of moles (or number of molecules) of component  $i$  [4]. Alternatively,  $X$  is defined as the composition-weighted sum of partial molar properties of the constituent components in the system:  $X = \sum_i n_i \bar{x}_i$  [5]. Experimentally, obtaining partial molar properties at extreme conditions is difficult (*i.e.* high temperatures and pressures) [6–9]. Computation of partial molar properties using molecular simulation [5,10,11] is not straightforward, as partial molar properties cannot be determined as a function of atomic positions or momenta of a single configuration in a system [5,10,12–14]. To date, different approaches have been developed based on:

(1) direct numerical differentiation [1,10,15], (2) Widom's test particle insertion method [10,12,13], (3) Kirkwood-Buff integrals [16–21], and (4) expanded ensembles [12,22]. Recently, Josephson et al. computed partial molar properties by fitting extensive thermodynamic properties (*e.g.*, enthalpy or volume) as a function of the instantaneous number of molecules of each component [5]. For example, partial molar enthalpies in a binary system,  $\bar{h}_1$  and  $\bar{h}_2$ , can be obtained by fitting the equation  $H = n_1 \bar{h}_1 + n_2 \bar{h}_2$  to the instantaneous enthalpy  $H$  as a function of the number of molecules of components 1 and 2 ( $n_1$  and  $n_2$ , respectively). This requires an ensemble in which the number of molecules of each component fluctuates (*e.g.* the  $NPT$  version of the Gibbs ensemble [23,24], reaction ensemble [25–27], or grand-canonical ensemble [28,29]). The advantage of this method is that simulation or legacy data (*e.g.*,  $H(n_1, n_2)$ ) can be easily fitted to the linear regression model without additional requirements. The required simulation data files are small compared to trajectory files containing all particle positions [5]. The partial molar properties can be computed from these data files using a single line of code in software such as MATLAB [30].

\* Corresponding author.

E-mail address: [t.j.h.vlugt@tudelft.nl](mailto:t.j.h.vlugt@tudelft.nl) (T.J.H. Vlugt).



identical expressions for partial properties in the grand-canonical ensemble (constant volume). For two-component mixtures, analytic expressions for  $(\partial U/\partial n_i)_{V,T,n_{j\neq i}}$  are provided. In Section 3, the application is extended to systems in which the pressure is constant (instead of constant volume). In Section 4, the compositional dependence of partial molar properties is examined using non-linear regression and a Gaussian process model [44]. Our conclusions are summarized in Section 5.

**2. Fluctuations in the grand-canonical ensemble**

For multicomponent adsorption in the grand-canonical ensemble, the derivative  $(\partial U/\partial n_i)_{V,T,n_{j\neq i}}$  is obtained from [31]

$$\left(\frac{\partial U}{\partial n_i}\right)_{V,T,n_{j\neq i}} = \sum_{i=1}^k \left(\frac{\partial U}{\partial(\beta\mu_k)}\right)_{V,T,\mu_{j\neq k}} \left(\frac{\partial(\beta\mu_k)}{\partial n_i}\right)_{V,T,n_{j\neq i}} \tag{6}$$

where  $k$  is the number of components. The second term on the right hand side of Eq. (6) is the change in chemical potential of component  $i$  while the number of molecules of all other components are fixed. Using a matrix notation, Eq. (6) can be written as

$$\begin{bmatrix} \left(\frac{\partial U}{\partial n_1}\right)_{n_{j\neq 1}} \\ \left(\frac{\partial U}{\partial n_2}\right)_{n_{j\neq 2}} \\ \vdots \\ \left(\frac{\partial U}{\partial n_k}\right)_{n_{j\neq k}} \end{bmatrix} = \underbrace{\begin{bmatrix} \left(\frac{\partial(\beta\mu_1)}{\partial n_1}\right)_{n_{j\neq 1}} & \left(\frac{\partial(\beta\mu_1)}{\partial n_2}\right)_{n_{j\neq 2}} & \cdots & \left(\frac{\partial(\beta\mu_1)}{\partial n_k}\right)_{n_{j\neq k}} \\ \left(\frac{\partial(\beta\mu_2)}{\partial n_1}\right)_{n_{j\neq 1}} & \left(\frac{\partial(\beta\mu_2)}{\partial n_2}\right)_{n_{j\neq 2}} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \left(\frac{\partial(\beta\mu_k)}{\partial n_1}\right)_{n_{j\neq 1}} & \cdots & \cdots & \left(\frac{\partial(\beta\mu_k)}{\partial n_k}\right)_{n_{j\neq k}} \end{bmatrix}}_{\mathbf{M}} \begin{bmatrix} \left(\frac{\partial U}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} \\ \left(\frac{\partial U}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} \\ \vdots \\ \left(\frac{\partial U}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \end{bmatrix} \tag{7}$$

in which we have defined matrix  $\mathbf{M}$ . The constant temperature and volume notation is dropped to make the equations more compact. In the grand-canonical ensemble, the derivatives  $(\partial U/\partial(\beta\mu_i))_{\mu_{j\neq i}}$  in Eq. (7) can be expressed as the covariance between  $U$  and  $n_i$  [31,32,35]:

$$\left(\frac{\partial U}{\partial(\beta\mu_i)}\right)_{\mu_{j\neq i}} = f(U, n_i) \tag{8}$$

where  $f(U, n_i)$  equals  $\langle Un_i \rangle - \langle U \rangle \langle n_i \rangle$ . The terms  $f(U, n_i)$  can be calculated directly from simulations in the grand-canonical ensemble. However, the derivatives  $(\partial(\beta\mu_i)/\partial n_i)_{n_{j\neq i}}$  in the matrix  $\mathbf{M}$  (Eq. (7)) cannot be directly calculated as ensemble averages, as the chemical potential is imposed. The inverse of the matrix  $\mathbf{M}$  is obtained directly from the reciprocals of the derivatives in  $\mathbf{M}$  [45]. Refs. [31,45] show that

$$\begin{bmatrix} \left(\frac{\partial(\beta\mu_1)}{\partial n_1}\right)_{n_{j\neq 1}} & \left(\frac{\partial(\beta\mu_1)}{\partial n_2}\right)_{n_{j\neq 2}} & \cdots & \left(\frac{\partial(\beta\mu_1)}{\partial n_k}\right)_{n_{j\neq k}} \\ \left(\frac{\partial(\beta\mu_2)}{\partial n_1}\right)_{n_{j\neq 1}} & \left(\frac{\partial(\beta\mu_2)}{\partial n_2}\right)_{n_{j\neq 2}} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \left(\frac{\partial(\beta\mu_k)}{\partial n_1}\right)_{n_{j\neq 1}} & \cdots & \cdots & \left(\frac{\partial(\beta\mu_k)}{\partial n_k}\right)_{n_{j\neq k}} \end{bmatrix} \underbrace{\begin{bmatrix} \left(\frac{\partial n_1}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \left(\frac{\partial n_1}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} & \cdots & \left(\frac{\partial n_1}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \\ \left(\frac{\partial n_2}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \left(\frac{\partial n_2}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \left(\frac{\partial n_k}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \cdots & \cdots & \left(\frac{\partial n_k}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \end{bmatrix}}_{\mathbf{M}^{-1}} = \mathbf{I} \tag{9}$$

where  $\mathbf{I}$  is the identity matrix. The advantage of this formulation is that the elements of the matrix  $\mathbf{M}^{-1}$  can be calculated directly as a function of ensemble averages in the grand-canonical ensemble [31,45]:

$$\left(\frac{\partial n_i}{\partial(\beta\mu_i)}\right)_{\mu_{j\neq i}} = f(n_i, n_i) \tag{10}$$

where  $f(n_i, n_i)$  denotes the covariance  $\langle n_i n_i \rangle - \langle n_i \rangle \langle n_i \rangle$ . By multiplying both sides of Eq. (7) with the inverse matrix  $\mathbf{M}^{-1}$  we have

$$\begin{bmatrix} \left(\frac{\partial n_1}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \left(\frac{\partial n_1}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} & \cdots & \left(\frac{\partial n_1}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \\ \left(\frac{\partial n_2}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \left(\frac{\partial n_2}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \left(\frac{\partial n_k}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} & \cdots & \cdots & \left(\frac{\partial n_k}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \end{bmatrix} \begin{bmatrix} \left(\frac{\partial U}{\partial n_1}\right)_{n_{j\neq 1}} \\ \left(\frac{\partial U}{\partial n_2}\right)_{n_{j\neq 2}} \\ \vdots \\ \left(\frac{\partial U}{\partial n_k}\right)_{n_{j\neq k}} \end{bmatrix} = \begin{bmatrix} \left(\frac{\partial U}{\partial(\beta\mu_1)}\right)_{\mu_{j\neq 1}} \\ \left(\frac{\partial U}{\partial(\beta\mu_2)}\right)_{\mu_{j\neq 2}} \\ \vdots \\ \left(\frac{\partial U}{\partial(\beta\mu_k)}\right)_{\mu_{j\neq k}} \end{bmatrix} \tag{11}$$

Rewriting Eq. (11) using Eqs. (8) and (10) leads to

$$\begin{bmatrix} f(n_1, n_1) & f(n_2, n_1) & \cdots & f(n_k, n_1) \\ f(n_1, n_2) & f(n_2, n_2) & \cdots & f(n_k, n_2) \\ \vdots & \vdots & \ddots & \vdots \\ f(n_1, n_k) & f(n_2, n_k) & \cdots & f(n_k, n_k) \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_k \end{bmatrix} = \begin{bmatrix} f(U, n_1) \\ f(U, n_2) \\ \vdots \\ f(U, n_k) \end{bmatrix} \tag{12}$$

To make the equations more compact, the terms  $d_i \in [1,k]$  are used to denote the derivatives  $(\partial U/\partial n_i)_{n_{j\neq i}}$ . We can show that Eq. (12) is identical to the set of equations for multiple linear regression (Eq. (4)) [36]. Writing out the  $l^{\text{th}}$  row of Eq. (12) leads to

$$d_1 f(n_1, n_l) + d_2 f(n_2, n_l) \dots + d_k f(n_k, n_l) = f(U, n_l) \tag{13}$$

Combining Eqs. (13) and (8), (10) leads to

$$\langle U \rangle \langle n_l \rangle - [d_1 \langle n_1 \rangle \langle n_l \rangle + d_2 \langle n_2 \rangle \langle n_l \rangle \dots + d_l \langle n_l \rangle \langle n_l \rangle \dots + d_k \langle n_k \rangle \langle n_l \rangle] + [d_1 \langle n_1 n_l \rangle + d_2 \langle n_2 n_l \rangle \dots + d_l \langle n_l n_l \rangle \dots + d_k \langle n_k n_l \rangle] = \langle U n_k \rangle \tag{14}$$

which can be rearranged to

$$\langle n_l \rangle \left[ \underbrace{\langle U \rangle - (d_1 \langle n_1 \rangle + d_2 \langle n_2 \rangle \dots + d_k \langle n_k \rangle)}_{d_0} \right] + d_1 \langle n_1 n_l \rangle + d_2 \langle n_2 n_l \rangle \dots + d_l \langle n_l n_l \rangle \dots + d_k \langle n_k n_l \rangle = \langle U n_k \rangle \tag{15}$$

In Eq. (15), we have defined the term  $d_0$ . For a sample size of  $N_s$ , this leads to the following independent linear equation

$$N_s d_0 + d_1 \sum_{i=1}^{N_s} n_{1i} + d_2 \sum_{i=1}^{N_s} n_{2i} \dots + d_k \sum_{i=1}^{N_s} n_{ki} = \sum_{i=1}^{N_s} U_i \quad (16)$$

Multiplying Eq. (15) with the sample size  $N_s$  and rearranging leads to

$$d_0 \sum_{i=1}^{N_s} n_{li} + d_1 \sum_{i=1}^{N_s} n_{li} n_{1i} + d_2 \sum_{i=1}^{N_s} n_{li} n_{2i} \dots + d_l \sum_{i=1}^{N_s} n_{li}^2 \dots + d_k \sum_{i=1}^{N_s} n_{li} n_{ki} = \sum_{i=1}^{N_s} U_i n_{li} \quad (17)$$

Eqs. (16) and (17) are identical to the equations for least squares multiple linear regression as in Eq. (4). Therefore, it is shown that identical expressions are obtained for (1) thermodynamic derivatives using least square multiple linear regression, and (2) thermodynamic fluctuations in the grand-canonical ensemble. While we demonstrate this equivalence for calculation of  $(\partial U / \partial n_i)_{V,T,n_{j \neq i}}$ , any extensive property  $X$  may be substituted for  $U$ .

For a two component system, it is still possible to solve Eq. (6) without resorting to matrix algebra. Using multiple linear regression ( $U = b_0 + b_1 n_1 + b_2 n_2 + \epsilon$ ), the partial derivatives in the grand-canonical ensemble can be written as

$$\left( \frac{\partial U}{\partial n_1} \right)_{V,T,n_2} = \frac{(\langle U n_1 \rangle - \langle U \rangle \langle n_1 \rangle)(\langle n_2^2 \rangle - \langle n_2 \rangle^2) - (\langle U n_2 \rangle - \langle U \rangle \langle n_2 \rangle)(\langle n_1 n_2 \rangle - \langle n_1 \rangle \langle n_2 \rangle)}{(\langle n_1^2 \rangle - \langle n_1 \rangle^2)(\langle n_2^2 \rangle - \langle n_2 \rangle^2) - (\langle n_1 n_2 \rangle - \langle n_1 \rangle \langle n_2 \rangle)^2} \quad (18)$$

and

$$\left( \frac{\partial U}{\partial n_2} \right)_{V,T,n_1} = \frac{(\langle U n_2 \rangle - \langle U \rangle \langle n_2 \rangle)(\langle n_1^2 \rangle - \langle n_1 \rangle^2) - (\langle U n_1 \rangle - \langle U \rangle \langle n_1 \rangle)(\langle n_1 n_2 \rangle - \langle n_1 \rangle \langle n_2 \rangle)}{(\langle n_1^2 \rangle - \langle n_1 \rangle^2)(\langle n_2^2 \rangle - \langle n_2 \rangle^2) - (\langle n_1 n_2 \rangle - \langle n_1 \rangle \langle n_2 \rangle)^2} \quad (19)$$

Similar to Eq. (5), the intercept is obtained from

$$b_0 = \langle U \rangle - b_1 \langle n_1 \rangle - b_2 \langle n_2 \rangle \quad (20)$$

It is important to note that computing fluctuations in the number of molecules in the grand-canonical ensemble is not always computationally efficient [11]. This is for example the case at high molecule densities (where insertion and removal of molecules is difficult [11]), or near inflection points in the adsorption isotherm (where  $(\partial n / \partial \mu) \approx 0$ ) [38–40]. In those cases, one may consider linear regression on a collection of  $NVT$  ensembles with varying number of molecules, which is similar to the numerical differentiation approach [1,10,15]. We will not consider this further here.

### 3. Application to constant pressure ensembles

In the previous sections, we have considered systems at constant volume. Partial molar properties are defined as partial derivatives of thermodynamic properties with respect to the number of molecules of component  $i$  at constant pressure, while keeping the number of molecules of all other components constant [1,10]. To calculate partial molar properties, one can either perform simulations directly at constant pressure, or perform a translation from a constant volume ensemble to a constant pressure ensemble as explained in Ref. [16]. For convenience, we choose to perform our simulations in constant pressure ensembles [46,47], such as the reactive isothermal-isobaric ensemble [25–27], or the isothermal-isobaric version of the Gibbs ensemble [24]. It is not practical to perform simulations in which only the pressure, temperature and the chemical potential are fixed ( $\mu PT$  ensemble) [11], as the pressure, temperature, and chemical potential are all intensive variables. In this ensemble, the size of the system in the simulation may decrease or increase without bound, making simulations generally unstable [48]. To impose an upper bound on the system

size, at least one extensive variable should be fixed (e.g. in the grand-canonical ensemble, the volume is fixed). To constrain the system size at constant pressure, there are several possibilities: (1) Performing simulation in the  $NPT$  version of the Gibbs Ensemble (where the total number of molecules is fixed) [11,24]; (2) Performing simulations in the  $NPT$  version of the reaction ensemble (total number of atoms is fixed) [25,27,49–51]. The reaction ensemble can also be considered as a grand-canonical ensemble in which the chemical potentials of reactants and reaction products are imposed in such a way that chemical equilibrium is obtained [25–27]. To calculate the partial derivatives of thermodynamic properties at constant pressure, one can fit the multiple linear regression model (Eq. (4)) to the simulation data [5]. Since the partial derivatives in Eqs. (8) and (10) do not change when the pressure is kept constant instead of the volume, the resulting expressions from fluctuations (Eq. (6)) do not change either. At constant pressure, it is possible to show that the intercept from the linear regression model will be zero [5]. As an example, one can write the instantaneous total volume as the sum of the composition-weighted partial molar volumes ( $k$  components) [5]

$$V = n_1 \bar{v}_1 + n_2 \bar{v}_2 \dots + n_k \bar{v}_k \quad (21)$$

in which  $\bar{v}_i$  is the partial molar volume of component  $i$ . By fit-

ting the least squares multiple linear regression model, ( $V = n_0 + n_1 \bar{v}_1 + n_2 \bar{v}_2 + \dots + n_k \bar{v}_k + \epsilon$ , in which  $\epsilon$  represents noise), the expression for the intercept follows from 5:

$$n_0 = \langle V \rangle - \bar{v}_1 \langle n_1 \rangle - \bar{v}_2 \langle n_2 \rangle - \dots - \bar{v}_k \langle n_k \rangle \quad (22)$$

From Eqs. (21) and (22) it follows directly that the intercept  $n_0$  should be zero, and this can be used as a test for the numerical calculations. We have carried out such a test for the systems in this study at constant pressure and observed that the values obtained for the intercept are statistically indistinguishable from zero. Nonetheless, we recommend constraining the intercept to be zero when capturing partial molar properties.

When applied to the reaction ensemble [26,27,49,50,52–55], the regression method proposed by Josephson et al. [5] can be used to calculate the reaction enthalpy without computing partial molar enthalpies of each species in the mixture. As an example, we consider the ammonia synthesis reaction ( $N_2 + 3H_2 \rightleftharpoons 2NH_3$ ) from our earlier work [12]. To facilitate reaction trial moves, the Continuous Fractional Component (CFC) version of the reaction ensemble was used [49]. The reaction enthalpy of the ammonia synthesis reaction is obtained directly by fitting the multiple linear regression model to instantaneous data obtained from a single simulation. Here, enthalpy of the system ( $H$ ) and the number of nitrogen molecules ( $n_{N_2}$ ) are related by

$$H = b_0 + b_1 n_{N_2} + \epsilon \quad (23)$$

in which  $b_1$  and  $b_0$  are the slope and intercept of the fitted regression line, respectively, and  $\epsilon$  is the error term. Just as in the grand-canonical ensemble, the intercept  $b_0$  will not be zero here. From simple linear regression (Eq. (1)) it follows that

$$b_1 = \frac{\langle H n_{N_2} \rangle_{\text{R}x\text{MC}} - \langle H \rangle_{\text{R}x\text{MC}} \langle n_{N_2} \rangle_{\text{R}x\text{MC}}}{\langle n_{N_2}^2 \rangle_{\text{R}x\text{MC}} - \langle n_{N_2} \rangle_{\text{R}x\text{MC}}^2} \quad (24)$$



**Table 1**

Residual reaction enthalpy ( $\Delta\bar{h}$ ) of the ammonia synthesis reaction, per mole of  $N_2$ , at  $T = 573$  K and  $P = 400$  bar and  $P = 800$  bar, computed both from simulations in the reaction ensemble (RxMC) [49] using Eq. (24), and the CFC-NPT ensemble [12,22]. The reported residual reaction enthalpies are relative to the reaction enthalpy at ideal gas conditions, as the contributions from the enthalpies of formation of  $N_2$ ,  $H_2$ , and  $NH_3$  are not included here in  $\Delta\bar{h}$ . The magnitude of the fluctuations for the number of ammonia molecules,  $\xi_{NH_3} = \sqrt{\langle n_{NH_3}^2 \rangle - \langle n_{NH_3} \rangle^2} / \langle n_{NH_3} \rangle$ , in the reaction ensemble is provided. Numbers in brackets are standard deviations of average volumes ( $\langle V \rangle$ ), from 5 independent simulations.

P/[bar]	RxMC [49]		$\langle V \rangle / \text{\AA}^3$	CFC-NPT [12]
	$\Delta\bar{h} / [\text{kJ mol}^{-1}]$	$\xi_{NH_3}$		$\Delta\bar{h} / [\text{kJ mol}^{-1}]$
400	-45.6(6)	0.0120(3)	81,800(200)	-45(2)
800	-75(1)	0.0070(1)	40,000(100)	-75(2)

in which the brackets  $\langle \dots \rangle_{\text{RxMC}}$  denote an ensemble average in the reaction ensemble. The slope in Eq. (24) is the reaction enthalpy per mole of nitrogen ( $b_1 = \Delta\bar{h}$ ). Similarly, the reaction enthalpy can be calculated by fitting the total enthalpy of the system as a function of the number of ammonia or hydrogen molecules. Since the fluctuations of both reactants and reaction products depend on the extent of the reaction, the final results will be identical. To verify Eq. (24), simulations of the ammonia synthesis reaction at  $P = 400$  bar and  $P = 800$  bar are performed in the CFC version of the reaction ensemble at  $T = 573$  K in a similar manner as described in Ref. [49]. The results are compared to the reaction enthalpies obtained from multiple simulations in the Continuous Fractional Component NPT (CFC-NPT) ensemble as in Ref. [12]. In these simulations, the reaction enthalpy follows from the partial molar enthalpies of all reactants and reaction products in the system, which are computed separately. For simulation details, the reader is referred to Refs. [12,49]. The resulting residual reaction enthalpies are shown in Table 1. The residual  $\Delta\bar{h}$  reported in Table 1 is with respect to the ideal gas contribution (102.07 kJ per mole of  $N_2$  [12]). At  $P = 400$  bar and  $P = 800$  bar, the residual  $\Delta\bar{h}$  values are non-zero which shows non-ideal behavior of the system at high pressures. Excellent agreement is observed between the results from the CFC-NPT ensemble simulations and direct computation of the reaction enthalpy using least squares linear regression. As shown in Table 1, the relative magnitude of the fluctuations in the reaction ensemble, encountered throughout a single simulation is small. Therefore, partial molar properties do not change in these simulations. It is important to note that for simulations both in the reaction and CFC-NPT ensembles, the enthalpy of the system includes contributions from so-called fractional molecules [35]. When the number of fractional molecules is less than 1% of the total number of molecules (as is the case here), the presence of fractional molecules does not significantly (given current computational resources) affect the computed ensemble averages [56].

#### 4. Compositional dependence of partial molar properties

The simplistic linear regression model Eq. (3) assumes constant partial molar properties. However, this assumption can be removed to allow for more flexibility in the choice of the model [5]. Non-linear models can detect and characterize non-linearities in the system. A non-linear model that still yields approximately constant partial molar properties would further strengthen the justification that the linear model is correct for the subset of compositions encountered in simulations of a system with more than 100 molecules for the minority component. Of course, for the global composition space, compositional dependence of partial molar properties cannot be neglected.

More generally, a function  $Y = F(\mathbf{n})$  predicting an extensive thermodynamic property  $Y$  from the numbers of molecules in the system can be regressed using the trajectory of a simulation, starting from an equilibrated state. When  $F(\mathbf{n})$  is a non-linear function, the partial molar properties obtained from regressing the trajectory will be dependent on the composition of the system. Therefore, a more robust alternative compared to linear regression, for example, Gaussian processes [44,57], may be employed as the functional form of  $F(\mathbf{n})$  in a purely data-driven perspective. It is important to note that  $F(\mathbf{n})$  cannot be arbitrary as  $Y$  is an extensive property. When  $F(\mathbf{n})$  is a Gaussian process, the predicted extensive property may not double when the number of molecules in the system is doubled, making the model thermodynamically inconsistent. Therefore, the fundamental mathematical feature of an extensive property should be considered when constructing the functional form of  $F(\mathbf{n})$ , that is, it needs to be homogeneous of  $\mathbf{n}$  with degree 1, that is,

$$\alpha F(\mathbf{n}) = F(\alpha \mathbf{n}) \quad (25)$$

Although  $F(\mathbf{n})$  is thermodynamically constrained, an arbitrary differentiable function operating on the mole fractions of the system  $f(\mathbf{x})$  can be used to satisfy this constraint, that is,

$$Y = F(\mathbf{n}) = Nf\left(\frac{\mathbf{n}}{N}\right) = Nf(\mathbf{x}) \quad (26)$$

where  $N = \sum_{k=1}^k n_k$ , and  $N$  is the total number of molecules in the system for each frame in the simulation trajectory. The non-linearity in  $f(\mathbf{x})$  allows for composition dependence of partial molar properties, and  $f(\mathbf{x})$  can be any type of model including polynomials, neural networks [58], or Gaussian processes [44,57,59]. Since  $f(\mathbf{x}) = Y/N$ , it is fitted by regressing the molar property as a function of the mole fractions.

The partial molar properties can be analytically evaluated given the partial derivatives (gradient) of  $f(\mathbf{x})$ . It is important to note that since mole fractions of  $n$  components sum to 1, the partial molar volumes cannot be simply evaluated as differentiating  $f(\mathbf{x})$  with respect to molar fractions [16]. In Ref. [1], it is shown that the partial molar property of  $Y$  for the  $i$ th component (assuming constant  $T$  and  $P$ ) can be expressed as

$$y_i = \left( \frac{\partial Y}{\partial n_i} \right)_{P,T,n_{j \neq i}} = \left( \frac{\partial (Nf(\mathbf{x}))}{\partial n_i} \right)_{P,T,n_{j \neq i}} \quad (27)$$

$$= f(\mathbf{x}) + \left( \frac{\partial f(\mathbf{x})}{\partial x_i} \right)_{P,T,x_{j \neq i}} - \sum_{j=1}^k x_j \left( \frac{\partial f(\mathbf{x})}{\partial x_j} \right)_{P,T,x_{j \neq i}}$$

For a linear function of  $f(\mathbf{x})$ ,  $y_i$  is a constant number because the second term and the third term always cancel. When a linear combination of mole fractions is used for  $f(\mathbf{x})$ , the well-known approach of calculating partial molar properties from Eqs. (27) and (26) is equivalent to multiple linear regression (Eq. (21)). Although the regression data for Eq. (26) are the mole fractions instead of the numbers of molecules, the effect of the fluctuating total number of molecules is negligible, and Eqs. (26) and (21) yield statistically indistinguishable partial molar properties.

As a comparison with linear regression, the partial molar properties for a natural gas condensate system containing methane,  $n$ -butane, and  $n$ -decane [5] were also calculated using a quadratic function using

$$f(\mathbf{x}) = \mathbf{x}C_2\mathbf{x}^T + \mathbf{x}C_1 + C_0 \quad (28)$$

and a Gaussian process [44,57] using

$$f(\mathbf{x}) = GP(\mathbf{x}) \quad (29)$$

Both models take the vector of mole fractions as the input and outputs the extensive property  $Y$  divided by the total number of

**Table 2**

Partial molar properties (including intramolecular potential energy contributions but excluding kinetic energy terms) from regressing simulation trajectories of natural gas condensates. Since the conformational distributions of large, flexible molecules may differ between liquid and vapor phases (*i.e.*, good and poor solvents) [60], the internal potential energy and enthalpy need to include the intramolecular potential energy contributions. The thermodynamic constraints for the Gibbs ensemble of the ternary mixture were  $N_{C1} = 1276$ ,  $N_{C4} = 425$ , and  $N_{C10} = 125$ ,  $P = 16.22$  MPa, and  $T = 333$  K. In the two-box Gibbs ensemble, the number fluctuations in the two boxes are identical, and the relative fluctuations are:  $\xi_{C1} = 0.109_4$ ,  $\xi_{C4} = 0.065_2$ , and  $\xi_{C10} = 0.034_1$  in the liquid phase and  $\xi_{C1} = 0.092_4$ ,  $\xi_{C4} = 0.176_5$ , and  $\xi_{C10} = 0.50_1$  in the gas phase. The mean values and uncertainties were obtained from 64 independent simulations. The input data for regression contains instantaneous values of the relevant properties from for each independent simulation at 100-cycle intervals. Parameters for quadratic regression were obtained using the analytical form of least-squares regression and hyperparameters for Gaussian process were selected using 5-fold cross-validation. For comparison, molar properties of the pure species are also provided for the stable phase at the same temperature and pressure as the ternary VLE simulations:  $N_{C1} = 1000$ ,  $N_{C4} = 800$ , or  $N_{C10} = 320$ , and  $P = 16.22$  MPa, and  $T = 333$  K in 16 independent NPT simulations (with 50,000 MC cycles each). Uncertainties are reported as the standard error of the mean from 64 or 16 independent simulations.

		Linear regression Sections 2 and 3	Quadratic regression Eq. (28)	Gaussian process Eq. (29)	Pure component
Liquid $\bar{V}_i$ [m <sup>3</sup> /mol] $\times 10^{-3}$	C1	0.0892 <sub>4</sub>	0.0891 <sub>4</sub>	0.0889 <sub>4</sub>	–
	C4	0.0927 <sub>9</sub>	0.0927 <sub>9</sub>	0.0929 <sub>9</sub>	0.102970 <sub>10</sub>
	C10	0.140 <sub>2</sub>	0.135 <sub>3</sub>	0.139 <sub>2</sub>	0.196578 <sub>14</sub>
MAD [nm <sup>3</sup> ]		1.26 <sub>1</sub>	1.26 <sub>1</sub>	1.26 <sub>2</sub>	
Vapor $\bar{V}_i$ [m <sup>3</sup> /mol] $\times 10^{-3}$	C1	0.1535 <sub>2</sub>	0.1520 <sub>5</sub>	0.1535 <sub>2</sub>	0.15221 <sub>3</sub>
	C4	–0.016 <sub>1</sub>	–0.015 <sub>1</sub>	–0.016 <sub>1</sub>	–
	C10	–0.253 <sub>6</sub>	–0.253 <sub>6</sub>	–0.265 <sub>7</sub>	–
MAD [nm <sup>3</sup> ]		3.23 <sub>1</sub>	3.22 <sub>1</sub>	3.24 <sub>4</sub>	
Liquid $\bar{U}_i$ [kJ/mol]	C1	–0.63 <sub>2</sub>	–0.63 <sub>2</sub>	–0.63 <sub>2</sub>	–
	C4	–11.38 <sub>7</sub>	–11.39 <sub>8</sub>	–11.40 <sub>8</sub>	–11.084 <sub>2</sub>
	C10	–17.8 <sub>1</sub>	–18.1 <sub>3</sub>	–17.8 <sub>1</sub>	–16.316
MAD [K]		1.712 <sub>8</sub> $\times 10^2$	1.710 <sub>8</sub> $\times 10^2$	1.721 <sub>1</sub> $\times 10^2$	
Vapor $\bar{U}_i$ [kJ/mol]	C1	–1.09 <sub>1</sub>	–1.12 <sub>1</sub>	–1.08 <sub>1</sub>	–1.5990 <sub>4</sub>
	C4	–8.26 <sub>5</sub>	–8.24 <sub>6</sub>	–8.26 <sub>6</sub>	–
	C10	–12.9 <sub>2</sub>	–12.9 <sub>2</sub>	–12.0 <sub>2</sub>	–
MAD [K]		94 <sub>1</sub>	94 <sub>1</sub>	94 <sub>1</sub>	
Liquid $\bar{H}_i$ [kJ/mol]	C1	0.79 <sub>2</sub>	0.78 <sub>3</sub>	0.78 <sub>3</sub>	–
	C4	–9.7 <sub>1</sub>	–9.8 <sub>1</sub>	–9.8 <sub>1</sub>	–9.420 <sub>4</sub>
	C10	–15.7 <sub>2</sub>	–16.0 <sub>4</sub>	–15.6 <sub>2</sub>	–13.134 <sub>14</sub>
MAD [K]		2.67 <sub>2</sub> $\times 10^2$	2.67 <sub>2</sub> $\times 10^2$	2.67 <sub>3</sub> $\times 10^2$	
Vapor $\bar{H}_i$ [kJ/mol]	C1	1.42 <sub>1</sub>	1.38 <sub>3</sub>	1.42 <sub>1</sub>	0.8716 <sub>13</sub>
	C4	–8.6 <sub>1</sub>	–8.6 <sub>1</sub>	–8.6 <sub>1</sub>	–
	C10	–16.8 <sub>3</sub>	–16.8 <sub>3</sub>	–16.2 <sub>4</sub>	–
MAD [K]		2.33 <sub>2</sub> $\times 10^2$	2.33 <sub>2</sub> $\times 10^2$	2.33 <sub>3</sub> $\times 10^2$	

molecules in the system. The parameters in Eq. (28) (quadratic model) were fitted using the analytic solution of polynomial least squares regression, and the functional form of Eq. (29) (Gaussian process) was calculated from the set of training data points in each cross-validation split [44]. The partial molar properties were calculated using Eq. (27), treating the mole fractions as independent variables. The code used for computing the partial molar properties can be downloaded from <https://github.com/SiepmannGroup/PartialMolarProperties>.

Table 2 lists the partial molar properties and mean absolute deviations (MAD) obtained by linear regression, quadratic regression, and a Gaussian process. In Table 2, MAD is the measure of how accurately the linear (or nonlinear) model captures the extensive property and it corresponds to the magnitude of the fluctuations in the system. For all cases except for the partial molar enthalpy in the vapor box where the Gaussian Process overfits each trajectory, all three models achieved almost perfect agreement in regressing the partial molar properties of the system. These results demonstrate that the quadratic model and Gaussian process have learned the same pattern as the linear model from the simulation trajectories. This shows the physical correctness of the linear regression method for the subset of composition space sampled by the simulations. It should be noted that the partial molar volumes for butane and decane in the gas phase are negative. Negative molar volumes are impossible for pure compounds at any state point, but can be observed for non-ideal mixtures. For comparison, we also carried out single-phase simulations for pure methane, butane, and decane in the isobaric-isothermal ensemble. The cho-

sen state point is above the critical point for pure methane and this phase is considered as vapor here, but butane and decane are found in stable liquid phases. Comparing the molar properties for pure systems to the partial molar properties in the corresponding phase (see Table 2) shows statistically different values for all molecules.

## 5. Conclusions

We have shown that in the grand-canonical ensemble, expressions for thermodynamic derivatives obtained from least squares multiple linear regression are identical to the expressions obtained from thermodynamic fluctuations. This provides a conceptually simple and computationally efficient approach to obtain thermodynamic properties from fluctuations in multicomponent systems. Multiple linear regression is thermodynamically consistent with fluctuations both in constant-volume and constant-pressure ensembles. We also show in the reaction ensemble that the reaction enthalpy can be obtained directly from a single simulation by fitting the enthalpy as a function of the number of reactant molecules with simple linear regression. In this work, we have assumed that composition range due to fluctuations encountered in a given simulation is small, so that partial molar properties are locally constant within a given simulation, noting that partial molar properties in non-ideal mixtures are not constant but vary over the global composition space. Nonlinear regression models capable of capturing compositional dependence of partial molar properties do not perform better than the linear model when applied to simu-

lations at a single state point, thus providing strong support that multiple linear regression captures the essential physics.

### Declaration of Competing Interest

None.

### Acknowledgments

This work was supported by NWO Exacte Wetenschappen (Physical Sciences) for the use of supercomputer facilities, with financial support from the [Nederlandse Organisatie voor Wetenschappelijk Onderzoek](#) (Netherlands Organization for Scientific Research, NWO). TJHV acknowledges NWO-CW for a VICI grant. This work was also supported by the Department of Energy, Office of Basic Energy Sciences, Division of Chemical Sciences, Geosciences and Biosciences, under Award DE-FG02-17ER16362 (TJR, YS, and JIS ). Computational resources from the Minnesota Supercomputing Institute are also gratefully acknowledged.

### References

- [1] S.I. Sandler, *Chemical, Biochemical, and Engineering Thermodynamics*, fourth ed., John Wiley & Sons, Hoboken, N.J., USA, 2006.
- [2] M.J. Moran, H.N. Shapiro, *Fundamentals of Engineering Thermodynamics*, fifth ed., John Wiley & Sons, West Sussex, England, 2006.
- [3] S.K. Schnell, P. Englebienne, J.-M. Simon, P. Krüger, S.P. Balaji, S. Kjelstrup, D. Bedeaux, A. Bardow, T.J.H. Vlugt, How to apply the Kirkwood–Buff theory to individual species in salt solutions, *Chem. Phys. Lett* 582 (2013) 154–157.
- [4] J.M. Smith, H.C. Van Ness, M.M. Abbott, *Introduction to Chemical Engineering Thermodynamics*, seventh ed., McGraw-Hill, New York, USA, 2005.
- [5] T.R. Josephson, R. Singh, M.S. Minkara, E.O. Fetisov, J.I. Siepmann, Partial molar properties from molecular simulation using multiple linear regression, *Mol. Phys.* 117 (2019) 3589–3602.
- [6] I.M. Abdulagatov, A.R. Bazaev, E.A. Bazaev, M.B. Saidakhmedova, A.E. Ramazanov, Volumetric properties of near-critical and supercritical water + pentane mixtures: molar, excess, partial, and apparent volumes, *J. Chem. Eng. Data* 43 (1998) 451–458.
- [7] I.M. Abdulagatov, A.R. Bazaev, E.A. Bazaev, S.P. Khokhlachev, M.B. Saidakhmedova, A.E. Ramazanov, Excess, partial, and molar volumes of n-alkanes in near-critical and supercritical water, *J. Solution Chem.* 27 (1998) 731–753.
- [8] H. Liu, J.P. O'Connell, On the measurement of solute partial molar volumes in near-critical fluids with supercritical fluid chromatography, *Ind. Eng. Chem. Res.* 37 (1998) 3323–3330.
- [9] T. Chang, R.W. Rousseau, P.K. Kilpatrick, Methanol synthesis reactions: calculations of equilibrium conversions using equations of state, *Ind. Eng. Chem. Process Des. Dev.* 25 (1986) 477–481.
- [10] P. Sindzingre, G. Ciccotti, C. Massobrio, D. Frenkel, Partial enthalpies and related quantities in mixtures from computer simulation, *Chem. Phys. Lett.* 136 (1987) 35–41.
- [11] D. Frenkel, B. Smit, *Understanding Molecular Simulation: From Algorithms to Applications*, second ed., Academic Press, San Diego, California, 2002.
- [12] A. Rahbari, R. Hens, I.K. Nikolaidis, A. Poursaeidesfahani, M. Ramdin, O.A. Moulτος, D. Dubbeldam, T.J.H. Vlugt, I.G. Economou, Computation of partial molar properties using continuous fractional component Monte Carlo, *Mol. Phys.* 116 (2018) 3331–3344.
- [13] P. Sindzingre, C. Massobrio, G. Ciccotti, D. Frenkel, Calculation of partial enthalpies of an argon–krypton mixture by NPT molecular dynamics, *Chem. Phys.* 129 (1989) 213–224.
- [14] B. Smit, D. Frenkel, Calculation of the chemical potential in the Gibbs ensemble, *Mol. Phys.* 68 (1989) 951–958.
- [15] S.M. Walas, *Phase Equilibria in Chemical Engineering*, Butterworth-Heinemann, USA, 1985.
- [16] S.K. Schnell, R. Skorpa, D. Bedeaux, S. Kjelstrup, T.J.H. Vlugt, J.M. Simon, Partial molar enthalpies and reaction enthalpies from equilibrium molecular dynamics simulation, *J. Chem. Phys.* 141 (2014) 144501.
- [17] J.G. Kirkwood, F.P. Buff, The statistical mechanical theory of solutions. I, *J. Chem. Phys.* 19 (1951) 774–777.
- [18] P. Krüger, S.K. Schnell, D. Bedeaux, S. Kjelstrup, T.J.H. Vlugt, J.-M. Simon, Kirkwood–Buff integrals for finite volumes, *J. Phys. Chem. Lett.* 4 (2012) 235–238.
- [19] N. Dawass, P. Krüger, J.-M. Simon, T.J.H. Vlugt, Kirkwood–Buff integrals of finite systems: shape effects, *Mol. Phys.* 116 (2018) 1573–1580.
- [20] N. Dawass, P. Krüger, S.K. Schnell, T.J.H. Vlugt, J.-M. Simon, Kirkwood–Buff integrals from molecular simulation, *Fluid Phase Equilib.* 486 (2019) 21–36.
- [21] P. Krüger, T.J.H. Vlugt, Size and shape dependence of finite-volume Kirkwood–Buff integrals, *Phys. Rev. E* 97 (2018) 051301–051305.
- [22] A. Rahbari, R. Hens, O.A. Moulτος, D. Dubbeldam, T.J.H. Vlugt, Multiple free energy calculations from single state point continuous fractional component Monte Carlo simulation using umbrella sampling, *J. Chem. Theory Comput.* 16 (2020) 1757–1767.
- [23] A.Z. Panagiotopoulos, Direct determination of phase coexistence properties of fluids by Monte Carlo simulation in a new ensemble, *Mol. Phys.* 61 (1987) 813–826.
- [24] A.Z. Panagiotopoulos, N. Quirke, M. Stapleton, D. Tildesley, Phase equilibria by simulation in the Gibbs ensemble: alternative derivation, generalization and application to mixture and membrane equilibria, *Mol. Phys.* 63 (1988) 527–545.
- [25] J.K. Johnson, A.Z. Panagiotopoulos, K.E. Gubbins, Reactive canonical Monte Carlo: a new simulation technique for reacting or associating fluids, *Mol. Phys.* 81 (1994) 717–733.
- [26] W.R. Smith, B. Triska, The reaction ensemble method for the computer simulation of chemical and phase equilibria. I. theory and basic examples, *J. Comput. Phys.* 100 (1994) 3019–3027.
- [27] H.C. Turner, J.K. Brennan, M. Lisal, W.R. Smith, K.J. Johnson, K.E. Gubbins, Simulation of chemical reaction equilibria by the reaction ensemble Monte Carlo method: a review, *Mol. Simul.* 34 (2008) 119–146.
- [28] D.J. Adams, Chemical potential of hard-sphere fluids by Monte Carlo methods, *Mol. Phys.* 28 (1974) 1241–1252.
- [29] G.E. Norman, V.S. Filinov, Investigations of phase transitions by a Monte-Carlo method, *High Temp. Res. USSR* 7 (1969) 216–222.
- [30] Optimization, *Toolbox User's Guide*, MathWorks, Inc, The MathWorks, USA, 2016.
- [31] F. Karavias, A.L. Myers, Isothermic heats of multicomponent adsorption: thermodynamics and computer simulations, *Langmuir* 7 (1991) 3118–3126.
- [32] S. Ban, T.J.H. Vlugt, F. Kapteijn, J.v den Bergh, Modeling the loading dependency of diffusion in zeolites: the relevant site model extended to mixtures in DDR-type zeolite, *J. Phys. Chem. C* 113 (2009) 21856–21865.
- [33] M.P. Allen, D.J. Tildesley, *Computer Simulation of Liquids*, second ed., Oxford University Press, Oxford, United Kingdom, 2017.
- [34] T.J.H. Vlugt, D. Dubbeldam, S. Ban, S. Calero, E. García-Pérez, Computing the heat of adsorption using molecular simulations: the effect of strong coulombic interactions, *J. Chem. Theory Comput.* 4 (2008) 1107–1118.
- [35] A. Torres-Knoop, A. Poursaeidesfahani, T.J.H. Vlugt, D. Dubbeldam, Behavior of the enthalpy of adsorption in nanoporous materials close to saturation conditions, *J. Chem. Theory Comput.* 13 (2017) 3326–3339.
- [36] R.E. Walpole, R.H. Myers, S.L. Myers, K. Ye, *Probability & Statistics for Engineers & Scientists*, ninth ed., Prentice Hall, Boston, USA, 2012.
- [37] A. Poursaeidesfahani, A. Torres-Knoop, M. Rigutto, N. Nair, D. Dubbeldam, T.J.H. Vlugt, Computation of the heat and entropy of adsorption in proximity of inflection points, *J. Phys. Chem. C* 120 (2016) 1727–1738.
- [38] T.J.H. Vlugt, R. Krishna, B. Smit, Molecular simulations of adsorption isotherms for linear and branched alkanes and their mixtures in silicalite, *J. Phys. Chem. B* 103 (1999) 1102–1118.
- [39] T.J.H. Vlugt, W. Zhu, F. Kapteijn, J.A. Moulijn, B. Smit, R. Krishna, Adsorption of linear and branched alkanes in the zeolite silicalite-1, *J. Amer. Chem. Soc.* 120 (1998) 5599–5600.
- [40] R. Krishna, B. Smit, T.J.H. Vlugt, Sorption-induced diffusion-selective separation of hydrocarbon isomers using silicalite, *J. Phys. Chem. A* 102 (1998) 7727–7730.
- [41] R.F. DeJaco, B. Elyassi, M.D. de, N. Mittal, M.M. Tsapatsis, J.I. Siepmann, Understanding the unique sorption of alkane- $\alpha$ ,  $\omega$ -diols in silicalite-1, *J. Comput. Phys.* 149 (2018) 072331.
- [42] I. Langmuir, The adsorption of gases on plane surfaces of glass, mica and platinum, *J. Amer. Chem. Soc.* 40 (9) (1918) 1361–1403.
- [43] M. Volmer, P. Mahnert, Über die auflösung fester körper in flüssigkeitsoberflächen und die eigenschaften der dabei entstehenden schichten, *Z. Phys. Chem.* 115 (1) (1925) 239–252.
- [44] M. Seeger, Gaussian processes for machine learning, *Int. J. Neural Syst.* 14 (2004) 69–106.
- [45] A.Z. Panagiotopoulos, R.C. Reid, On the relationship between pairwise fluctuations and thermodynamic derivatives, *J. Chem. Phys.* 85 (1986) 4650–4653.
- [46] W.W. Wood, Monte Carlo calculations for hard disks in the isothermal-isobaric ensemble, *J. Chem. Phys.* 48 (1968) 415–434.
- [47] W.W. Wood, NpT-ensemble Monte Carlo calculations for the hard-disk fluid, *J. Chem. Phys.* 52 (1970) 729–741, doi:10.1063/1.1673047.
- [48] J.R. Ray, Ensembles and computer simulation calculation of response functions, *Handb. Mater. Model.* (2005) 729–743.
- [49] A. Poursaeidesfahani, R. Hens, A. Rahbari, M. Ramdin, D. Dubbeldam, T.J.H. vlugt, Efficient application of continuous fractional component Monte Carlo in the reaction ensemble, *J. Chem. Theory Comput.* 13 (2017) 4452–4466.
- [50] T.W. Rosch, E.J. Maginn, Reaction ensemble Monte Carlo simulation of complex molecular systems, *J. Chem. Theory Comput.* 7 (2011) 269–279.
- [51] E.O. Fetisov, I.-F.W. Kuo, C. Knight, J. VandeVondele, T.V. Voorhis, J.I. Siepmann, First-principles Monte Carlo simulations of reaction equilibria in compressed vapors, *ACS. Cent. Sci.* 2 (2016) 409–415.
- [52] C.H. Turner, J.K. Johnson, K.E. Gubbins, Effect of confinement on chemical reaction equilibria: the reactions  $2\text{NO} \leftrightarrow (\text{NO})_2$  and  $\text{N}_2 + 3\text{H}_2 \leftrightarrow 2\text{NH}_3$  in carbon micropores, *J. Comput. Phys.* 114 (2001) 1851–1859.
- [53] S.P. Balaji, S. Gangarapu, M. Ramdin, A. Torres-Knoop, H. Zuilhof, E.L. Goetheer, D. Dubbeldam, T.J.H. Vlugt, Simulating the reactions of  $\text{CO}_2$  in aqueous monoethanolamine solution by reaction ensemble Monte Carlo using the continuous fractional component method, *J. Chem. Theory Comput.* 11 (2015) 2661–2669.
- [54] R.G. Mullen, S.A. Corcelli, E.J. Maginn, Reaction ensemble Monte Carlo simulations of  $\text{CO}_2$  absorption in the reactive ionic liquid triethyl(octyl)phosphonium 2-cyanopyrrolide, *J. Phys. Chem. Lett.* 9 (2018) 5213–5218.



- [55] M. Lísal, M. Bendová, W.R. Smith, Monte Carlo adiabatic simulation of equilibrium reacting systems: the ammonia synthesis reaction, *Fluid Phase Equilib.* 235 (2005) 50–57.
- [56] A. Rahbari, R. Hens, D. Dubbeldam, T.J.H. Vlugt, Improving the accuracy of computing chemical potentials in CFCMC simulations, *Mol. Phys.* 117 (2019) 3493–3508.
- [57] J.Q. Shi, T. Choi, *Gaussian Process Regression Analysis for Functional Data*, CRC Press, Boca Raton, USA, 2011.
- [58] S. Haykin, *Neural Networks: a Comprehensive Foundation*, Prentice Hall PTR, USA, 1994.
- [59] D.F. Specht, A general regression neural network, *IEEE Trans. Neural Netw.* 2 (1991) 568–576.
- [60] N.D. Zhuravlev, M.G. Martin, J.I. Siepmann, Vapor-liquid phase equilibria of triacontane isomers: deviations from the principle of corresponding states, *Fluid Phase Equilib.* 202 (2002) 307–324.