# データ駆動科学の立場からみた 物質科学と情報科学の接点

### 安藤 康伸

AIST, 機能材料コンピュテーショナルデザイン研究センター

# 自己紹介

- 名前:安藤康伸(36)
- 出身:愛知県西尾市

### <u>略歴</u>

- ✔ 2012.3. 博士(理学)@東大院理
- ✓ 2012.4~2013.4 ポスドク@AIST
- ✔ 2013.5~2016.3 助教 @東大マテエ
- ✓ 2016.4 ~ 研究員 @ 機能材料コンピューテショ
   ナルデザイン研究センター (CD-FMat), AIST

情報科学X物質科学

### **Back ground**

- ✔ 物性物理学
- ✔ 表面・界面科学
- ✔ 計算物質科学



### 電極一電解液界面の第一原理分子動力学





# 情報科学ブームの火付け役

## <u>Materials Genome Initiative (MGI from 2011)</u>

### <u>戦略目標</u>



- I. 材料イノベーションのインフラ整備
- 2. 先端材料開発に関する国家目標達成
- 3. 次世代の研究従事者育成

## データベース・機械学習の活用を掲げる

# 物質・材料科学???





# 製品の"原材料"を研究する仕事



## 情報的な視点から見た物質科学

## <u>ビッグデータというよりスモールデータ</u>

✔ 典型的な実験:高々100,ようやく1,000サンプル出せることも

✔ 他の実験室の結果とも直接比較が困難(環境・装置・人)



### <u>汎用的な物質表現・処理方法の不在</u>

✔ 対象:高分子・半導体・金属 etc.

✓ スケール:Inm ~ Im (9桁!!)



✔ 構造やシステムの表現方法の多様性

物質科学の強みである「制御性(再現性)」「理論(事前知 識)」を生かした中規模な情報処理スタイルの開発が必要

## 事例の蓄積が鍵

Appl. Phys. Lett. 91, 132102 (2007). Appl. Phys. Rev. 4, 011105 (2017).

ハイスループット計測データ解析

### High-Throughput Experimental (HTE) methodologies



FIG. 2. (a) Electrical conductivity, (b) Seebeck coefficient, and (c) power factor of the composition-spread  $(Ca_{1-x-y}Sr_xLa_y)_3Co_4O_9$  film (0 < x < 1/3) and 0 < y < 1/3). Reproduced with permission from Appl. Phys. Lett. **91**, 3 (2007). Copyright 2007 AIP Publishing LLC.<sup>56</sup>

## <u>Materials "Library"</u>:単一シートに複数の組成を連続的に散布 物質空間を一気に観測、最適な組成を高速に発見する

(観測に~6 h)

Y. N.-Ohno, M. Haze, Y. Yoshida, K. Hukushima, Y. Hasegawa, and M. Okada, J. Phys. Soc. Jpn. 85, 093702 (2016).

準粒子干渉測定のダウンサンプリング

### 課題: 『短い観測時間でも正しい結果を得たい』



**Fig. 1.** (a) dI/dV map of Ag(111) surface. (b) FT of (a) obtained by conventional method.

Sampling in R space can be reduced.

Chem. Lett. 47, 284–287 (2018).

計測データを計算データでフィット

#### 高活性ホスホン酸触媒をシミュレーションからデザインすることに成功(矢田・安藤・永田)



# 第一原理計算データを利用した 機械学習ポテンシャル研究

## 原子・分子シミュレーション

#### 安定構造探索





振動モード・反応座標

vinyl alcohol to acetaldehyde (NEB method)



DMol3 TS search (LST/QST) 51.473 cal/mol, -10.851 cal/mol (PBE)





LTO表面構造の同定

(Nat. Commun. 8, 15975 (2017))

- 静的な安定状態から動的な挙動まで
- 構造と特性・観測データを結びつけるキー
- ナノスケール物理では必須のツール

Phys. Chem. Chem. Phys., 2018, 20, 11586-11591

## NNポテンシャル発展の鍵

#### <u>表面化学反応のシミュレーション</u>

### <u>方法論に求められること</u>

1.電子状態変化が取り扱えること
 2.エネルギー超曲面上でのダイナミクス
 3.ダイナミクスの統計性



- ・電子状態変化を伴うため、DFT計算が必須
- ・DFTのみでは計算が大変で統計が集まらない

### NNによるモデリングの試みは

1995年,2004年にすでに報告あり

J. Chem. Phys. **103** (10), 8 (**1995**). Chem. Phys. Lett. 395, 210 (**2004**).



Pd表面でのH₂分解反応

## NNフィッティングの諸問題

### 表面構造の特徴量設計

fcc(111)の表面形状に合わせてフーリエ変換を元に構成する方法を提案[1]

[1] J. Behler, S. Lorenz, and K. Reuter, J. Chem. Phys. 127, 014705(2007).

### ただし問題に合わせて複雑な関数を設計することは決して 容易ではなくもっと簡便な表現が必要に。

### 3次元配置を入力にしたNNの問題

- 入力層が原子数に依存するため、異なる系には適用できない(汎用性)
- 同種粒子の入れ替えに対する対称性が破れてしまう
- 同様にポテンシャルが満たすべき対称性が保証されない

## 入力表現・汎用性・対称性がキーワード

Behler and M. Parrinello, Phys. Rev. Lett. 98, 146401 (2007).

## Behler - Parrinelloの方法



- 正味のINPUTは、サブネットの入力 x 原子数
- 学習すべきNNのパラメータは、原子種ごとに共通。
- また入力ベクトルGの次元は固定できるので原子数を増やすことは容易

# 機械学習ポテンシャル作成の構成技術



機械学習ポテンシャルは**(1)構造の記述(2)モデリング**の違いで大別 使用するメリット・目的で手法を使い分ける必要がある

## アモルファス中のイオン安定配置の全探索

### 大問題:高精度に計算できるDFTではコストがかかりすぎる

(N<sup>3</sup>回の構造最適化 x 計算時間 T) N~50, T~1 hour以上 > 10 year さらに拡散障壁の計算も必要...



### Behler-Parrinelloの方法でポテンシャル作成





# 焼きなまし法の冷却速度依存性



18

# 計算・計測データの自動解析技術

安藤, 藤掛, 渡邉, 表面科学 36, 515-520 (2015).

自明なマップ上の類似度による分類

### 課題:『たくさんあるデータを自動で分類したい』



## 他の「スペクトルが多すぎる」問題



白丸の上の数字が右のスペクトルデータの数字に対応 白丸が各データの得られた大体の位置に対応

二本のピークが綺麗に見える領域(上二つ)と見えない領域(下二つ) があった

#### (86 x 86 pixels) 7,396

Collaborated with 永村直佳(NIMS), 松村太郎次郎(PD), 永田賢二(NIMS), 赤穂昭太郎(AIST)

ピーク位置の自動推定

### 課題:『ピークの位置を大まかにでも自動抽出したい』





# 非線形最小二乗フィットはしんどい



### <u>混合ガウスモデル推定</u>

推定対象

- 1. 分布の平均 μ<sub>k</sub> と分散 σ<sup>2</sup>k, 混合率 π<sub>k</sub>
- 2. データ点n の潜在変数 r<sub>nk</sub>
- E (Expectation)-step: (1)をもとに(2)を推定

$$E[z_{nk}] = \frac{\pi_k N(x_n, \theta_k)}{\sum_j \pi_k N(x_n | \theta_k)} = \gamma(z_{nk}) \quad \longleftarrow \quad (1) \succeq x_n ic 依存$$

M (Maximization)-step:(2)をもとに(1)を推定

完全データ(データxnと潜在変数rnkの組)から対数尤度の期待値を計算して最大化

$$\mu_{k}^{\text{new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) x_{n}, \quad \sigma_{k}^{2 \text{ new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) (x_{n} - \mu_{k}^{\text{new}})^{2}, \quad \pi_{k}^{\text{new}} = \frac{N_{k}}{N}$$
繰り返すと尤度が単調増加!



図3-1 混合ガウス分布の例

 $N_k = \sum \gamma(z_{nk})$ 

n=1

実データへの適用上の課題

## 放射光スペクトルはイベント数が膨大(~10<sup>7</sup> events)

<u>混合ガウス分布のモデルと対数尤度</u>

## 放射光スペクトルのイベント数を適切に生かし かつ高速に動作するように改良が必要

# スペクトル解析に適したEMアルゴリズム 通常のEMアルゴリズムの入力: $X_{1D} = \{x_1, x_2, x_3, \dots, x_N\}$ イベント列 1 2 3 1 2 2 2 1 、 改良型EMアルゴリズムの入力: $X_{2D} = \begin{cases} \hat{x}_1, & \hat{x}_2, & \hat{x}_3, & \cdots, & \hat{x}_M \\ w_1, & w_2, & w_3, & \cdots, & w_M \end{cases}$ ヒストグラム形式 $w_1 = 3$ 1 2 2 2 改良型EMアルゴリズムの対数尤度 (ポイント: N >> M) $\ln p(\{x_1, \dots, x_N\}) = \sum_{k=1}^{N} \ln \left( \sum_{k=1}^{K} \pi_k N(x_n | \mu_k, \sigma_k) \right) \quad \forall x \in \mathbb{N} \setminus \mathbb{N} \setminus$ • ビン数はイベント総数より $= \sum_{n=1}^{M} w_n \ln \left( \sum_{k=1}^{K} \pi_k N(\hat{x}_n | \mu_k, \sigma_k) \right)$ E倒的に小さい ・ イベント数が増えても 定質同粉が恋せった

## 計算速度と精度に関する性能評価

#### 通常のEMアルゴリズム:計算負荷が測定イベント数に依存 🦛 改善 8



- ランダムな初期値でも安定して動作、極めて低コストでモデル推定(ピーク抽出)が可能
- 初期値を振って最良のモデルを探せばさらに精度をあげられる
- 高ノイズデータ・埋もれたピークに関してはモデリングが困難(事前知識が必要)

## やりたかったこと



白丸の上の数字が右のスペクトルデータの数字に対応 白丸が各データの得られた大体の位置に対応

#### 光学顕微鏡の観測結果



この辺りのピーク位置の違いは光学顕微鏡像に現れていない

### EMアルゴリズムによるピーク位置推定結果

(~12 h程度)

## まとめ:今後、物質科学者として目指したいこと

