

マテリアルズ・インフォマティクスの狙い

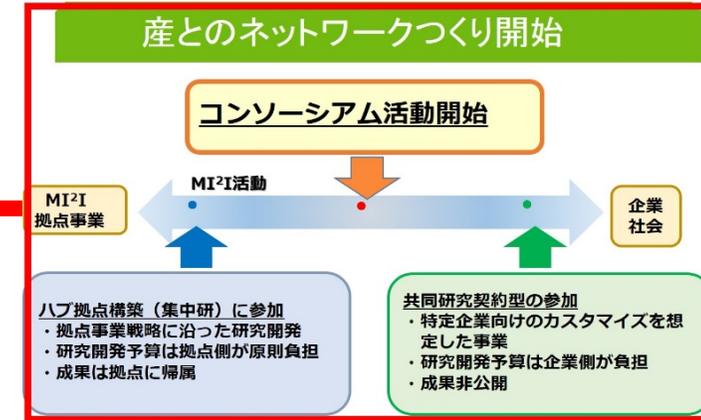
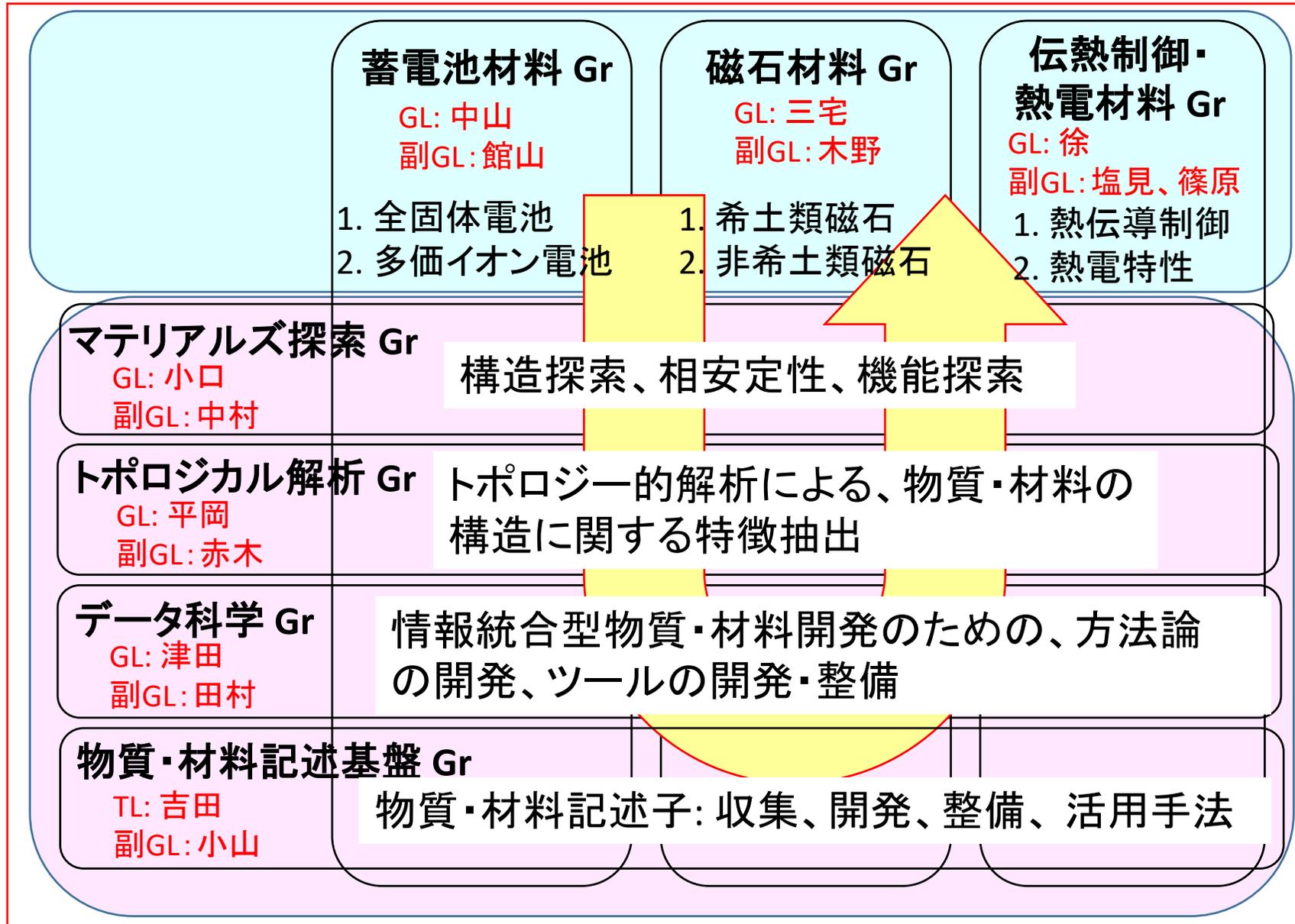
物質・材料研究機構 名誉フェロー・エグゼクティブアドバイザー
寺倉清之

マテリアルズ・インフォマティクスの現状と将来展望

2018/1/29

JSTイノベーションハブ構築支援事業 (H27-H31)

H29.9 / 51社
 コンソーシアム活動



再委託先



1. マテリアルズ・インフォマティクスとは何か？
2. データ科学のアプローチにおける
記述子と目的変数
3. ベイズ推定の意義と適用例
物理モデルをデータ解析に組み入れやすい。
データが少なくても、存在するデータを活かした推定ができる。
4. 「少ないデータの問題」への対応
転移学習 (+ multi-fidelity model)
5. 予測と理解

物質・材料開発 I

1957 BCS理論

その後、物性物理の終焉説が起こった。

その後の新物質発見の例

1976 カーボンナノチューブ(遠藤守信)

1985 フラレン

1986 銅酸化物超伝導体

1991 カーボンナノチューブ(飯島澄男)

2005 トポロジカル絶縁体の理論提案

2006 鉄系超伝導体

2007 トポロジカル絶縁体の実験的存在確認

これらの新物質発見が
基礎的な物性発現の機構
の提案と実証という
活発な物質科学の活動
の展開に繋がっている。

これらは、経験、洞察、ひらめき、偶然、などなどによる新奇物質発見

今のところ、データ科学的な研究とは明らかな繋がりをつけにくい。

- ・1/25 までの見解
- ・それ以後は、後述のように
少し変わった。

物質・材料開発 II

科学技術の進歩に後押しされた物質開発

この場合は、目標とする物性や機能が明確なものが多い。

物質設計

物質設計は、本来逆問題なので実現が困難

* Design principle によるスクリーニング

* データ科学の活用
機械学習 (bio-informatics, cheminformatics, materials informatics)

課題はSTAGEで異なる

真鍋明氏作成

	STAGE I 新物質創成 例) C ₆₀	STAGE II 物性極値化 Materials Genome	STAGE III 材料最適化 Integrated Computational Materials Engineering	STAGE IV 適用研究開発
内容	従来の特性限界超物質探索	結晶構造あり 元素置換 ドープ 極値を探す	材料化 プロセス・組織構造 の最適化	システム設計 試作実証 信頼性確保
ポイント	コンセプトひらめき 実験発見 Abduction	傾向予測と実験 Deduction Induction	実験検証 特性トレードオフ克服 Induction主体	Virtual Prototype シミュレーション Deduction
データ共有	Unknown OPEN 知識	OPEN 一部CLOSE 物質データベース	特性CLOSE/OPEN プロセスCLOSE	固有材料CLOSE 一般材料OPEN
MIの期待	逆問題 特性→構造予測	結晶構造・特性相関 QSPR	特性・組織相関 プロセス・組織相関	短期間化 (時間、費用)
課題	方法論研究	事例研究 手法選択	組織構造データ化 データ形式統一 (メタデータ)	各種シミュレーション

今後 ← 現状での主目的 → 今後

レシピ + 支配方程式

出口

計測インフォマティクス

計測の自動化
大規模実験施設 SPring8, J-PARC

STEM-EELS, Atom Probe Tomography,
SAXS, SANS

大型スパコン 京、ポスト京
不成功データの蓄積6

既にbig data の時代

- ・big data を活用した物質探索
- ・実験の効率化
- big data の自動解析

JST CREST, さきがけ
計測技術と高度情報処理の融合による
インテリジェント計測・解析手法の開発と応用

情報統合型研究の3つの主要目標

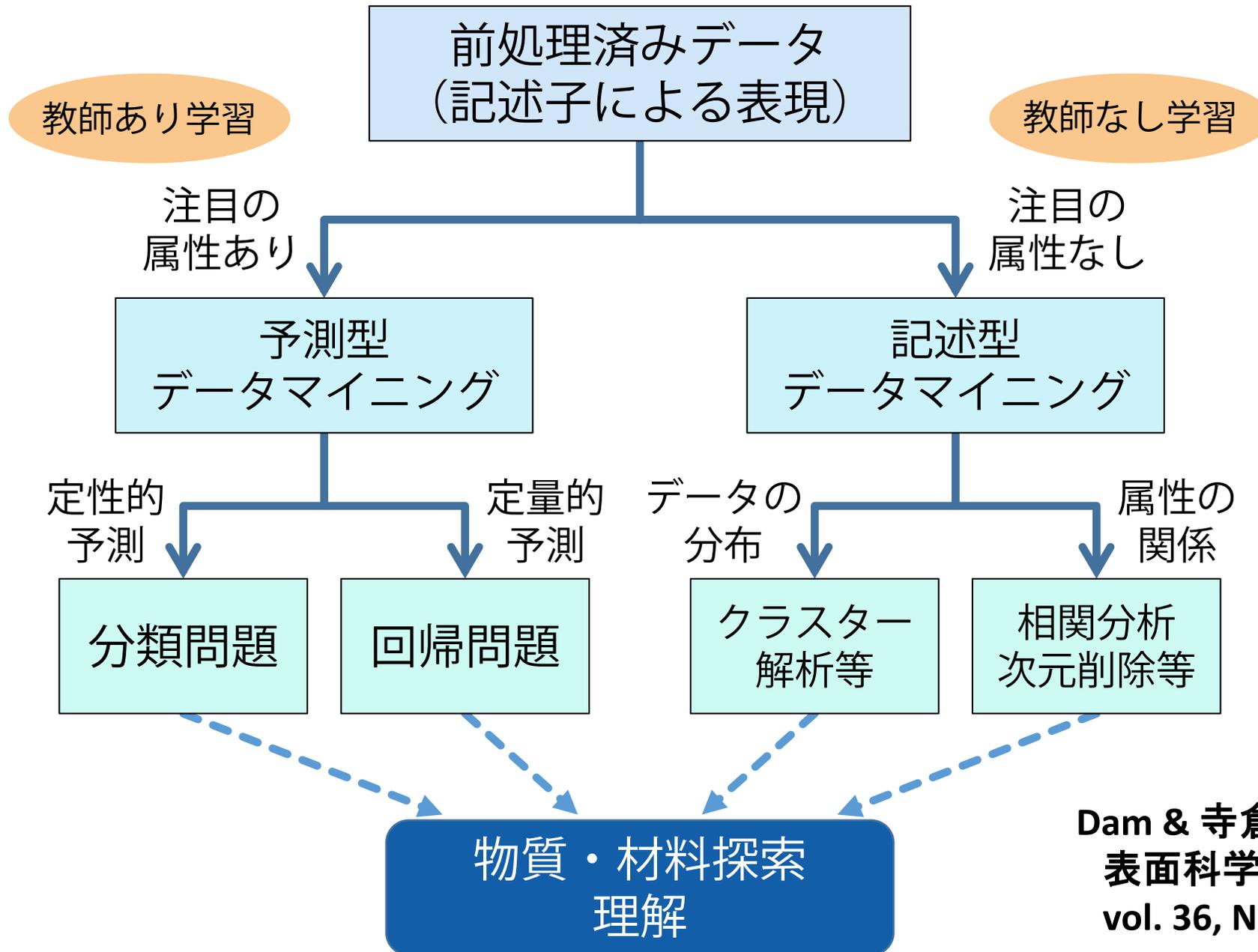
1. 解析から予測へ、予測から**設計**へ
解析型研究から**開拓型研究**へ

2. 複雑な現象や機能を制御している原因の**解明**

3. 実験・計算における解析型研究の**高速化**

蓄積したデータからの学習に基づく逆問題対応

個々の解析への対応
特に、計測インフォマティクス
ただし、これもデータ科学で上記の1, 2と同じ。



Dam & 寺倉
表面科学
vol. 36, No.10, p.507 (2015)

回帰解析

$$\text{機能: } F = f(x_1, x_2, \dots) \quad (1)$$

$$\text{記述子: } x_1, x_2, \dots \quad (2)$$

目的の機能 F と記述子 (x_1, x_2, \dots) の間の関係式 (1) を導くだけでなく、記述子 (x_1, x_2, \dots) を決める。

記述子: ①機能 F を制御する物理量であり、
②その機能に関して、対象の間の類似性を測るのに適している。

データからの回帰で得られる式 (1) は、**相関関係** と呼ばれる。
相関関係 は必ずしも**因果関係** ではない。 **➡ 予測の検証が必要**

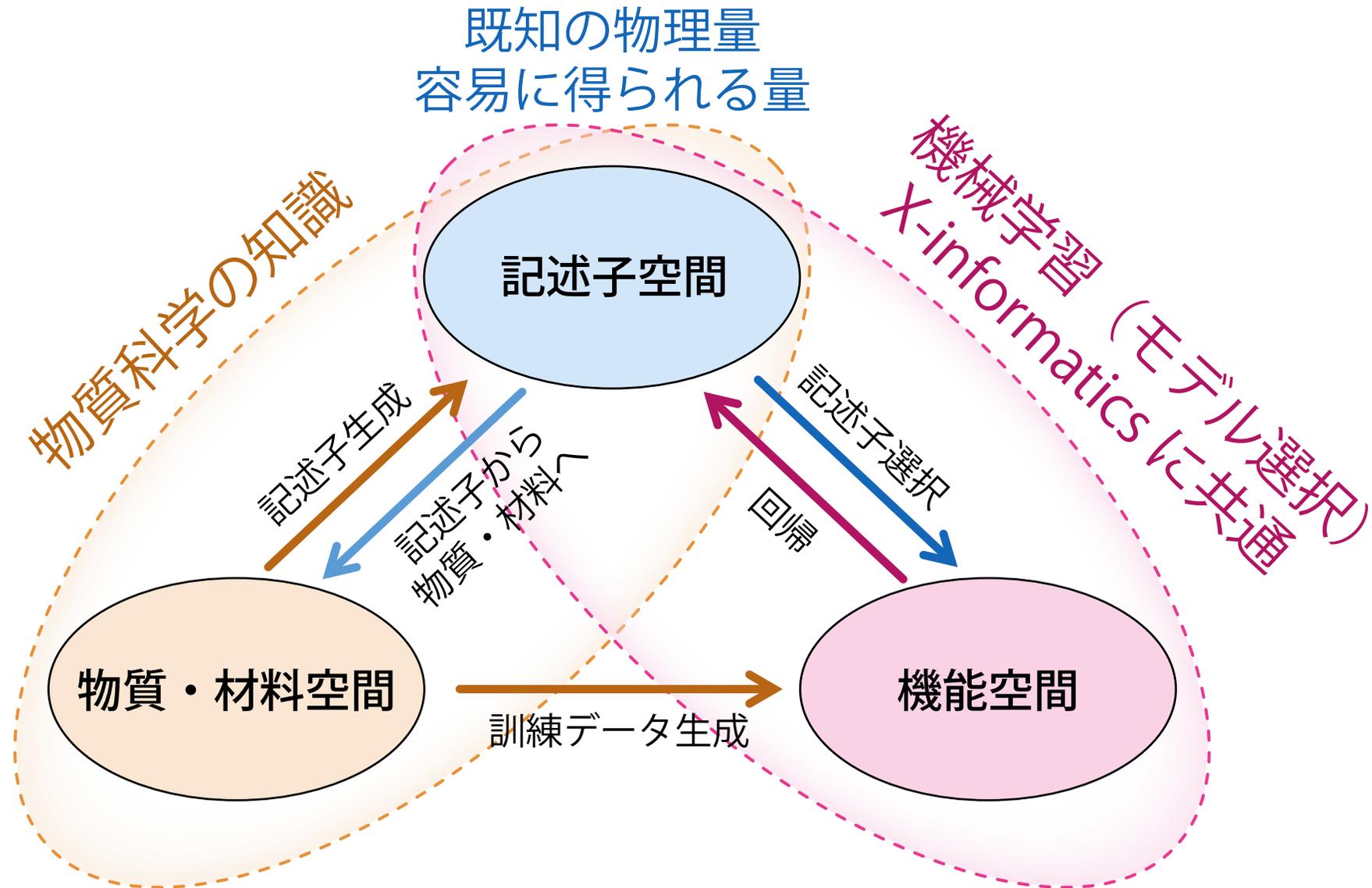
相関関係 が、**因果関係** になっているかどうかの判定もデータ科学の課題。

All models are wrong;
some models are useful.



George E. P. Box

物質・材料に機械学習を適用



現状では、
AIにはできない。

課題の設定：研究としては常にここが最重要

知性

教師あり学習の場合

目的変数（狙いの機能を測る量）の設定
測定や計算が容易な量が望ましい。

: 後述の転移学習に関係

うまく選ぶことが成功の鍵。専門知識が必要。

知能

記述子（目的変数を制御する量）の設定
既知、あるいは測定や計算が容易な量

記述子の整備は MI²¹でも重要課題
第5回チュートリアルセミナーで、吉田亮氏が説明

山本一成氏
(ポナンザ開発者)

機械学習の手法の選択

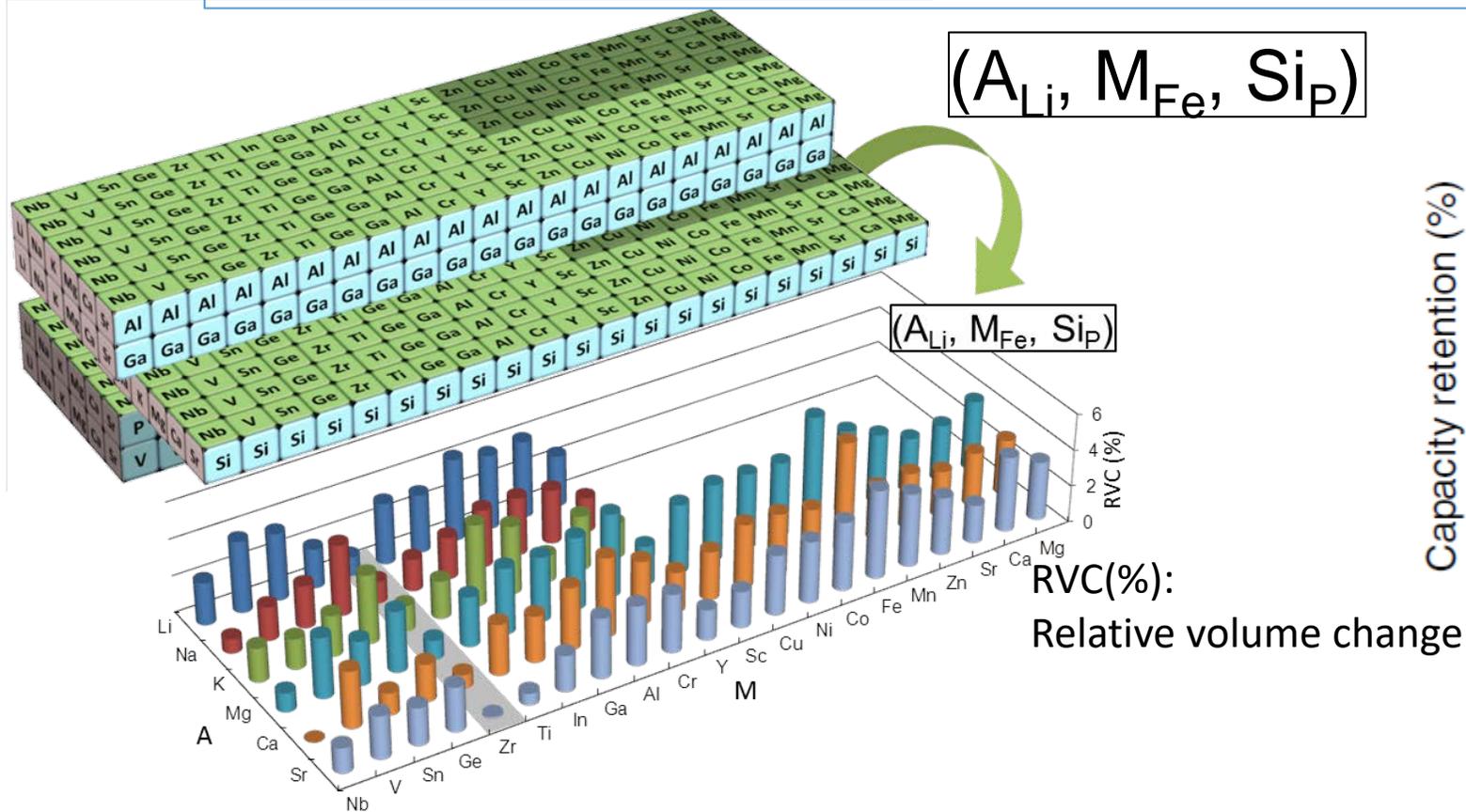
目的変数の選択の例

cosubstitution in LiFePO₄

京都大学
田中功 グループ

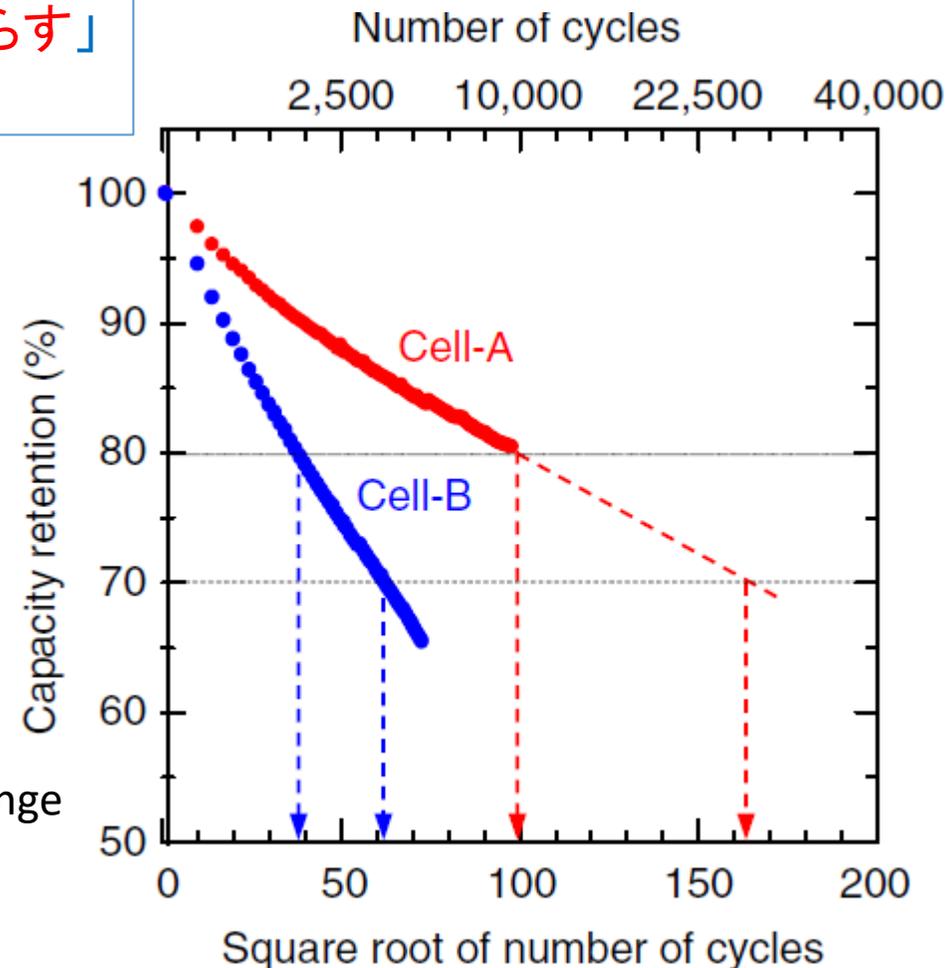
M. Nishijima et al. Nat. Commu. 5, 4553 (2014)

長寿命化を扱うのに、「充電時と放電時の体積変化を減らす」というデザインプリンシプルを立てた。



Total number of DFT calculations: 2,000

実証



寿命が6倍

ARTICLE

Received 2 Sep 2015 | Accepted 4 Mar 2016 | Published 15 Apr 2016

DOI: 10.1038/ncomms11241

OPEN

Accelerated search for materials with targeted properties by adaptive design

Dezhen Xue^{1,2}, Prasanna V. Balachandran¹, John Hogden³, James Theiler⁴, Deqing Xue² & Turab Lookman¹

the complex search space. Our strategy uses inference and global optimization to balance the trade-off between exploitation and exploration of the search space. We demonstrate this by finding very low thermal hysteresis (ΔT) NiTi-based shape memory alloys, with $\text{Ti}_{50.0}\text{Ni}_{46.7}\text{Cu}_{0.8}\text{Fe}_{2.3}\text{Pd}_{0.2}$ possessing the smallest ΔT (1.84 K). We synthesize and characterize 36 predicted compositions (9 feedback loops) from a potential space of $\sim 800,000$ compositions. Of these, 14 had smaller ΔT than any of the 22 in the original data set.

多元系形状記憶合金 Ti-Ni-Cu-Fe-Pd での温度ヒステリシス最小化
最適化組成を予測して、実証

実験とデータ科学の協働

9回の iteration で最適化実現

最適組成合金 $\text{Ti}_{50.0}\text{Ni}_{46.7}\text{Cu}_{0.8}\text{Fe}_{2.3}\text{Pd}_{0.2}$

機械学習の手法

1. Bayes 推定
 2. 転移学習
- に絞る。

Bayes 推定

同時確率: $P(X, Y) = P(Y | X)P(X) = P(X | Y)P(Y)$

$$\text{Bayes の定理: } P(X | Y) = \frac{\overset{\text{事後確率}}{P(Y | X)} \overset{\text{事前確率}}{P(X)}}{P(Y)} \propto \overset{\text{尤度}}{P(Y | X)} \overset{\text{事前確率}}{P(X)} \quad (1)$$

原因 結果
結果 原因
物質 物性
物性 物質

逆構造物性相関

構造物性相関

原因と結果の関係の向きを逆転できる。 → 逆問題に好都合。

与えられた結果を与える原因を求めるには、事後確率を最大にする x を探すことになる。

計算機の発達により、これは **MCMC (Markov-chain Monte Carlo) 法**などによって解かれるようになった。

観測データのセット $\{(y_i, \vec{x}_i); i=1 \sim N\}$ に対する解析

1. 尤度におけるモデルの活用 パラメータ推定

パラメータ群 Θ を含むモデル (法則) $G(\vec{x}_i; \Theta)$ により、

$$y_i = G(\vec{x}_i; \Theta) + \varepsilon_i$$

(5)

で表す。 ε_i はランダムノイズで、分散が σ^2 のガウス分布に

従うとして、尤度は

ノイズのバラつきと
データそのもののバラつきは別物

$$p(y_i | \vec{x}_i, \Theta) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right) \exp \left\{ -\frac{1}{2\sigma^2} (y_i - G(\vec{x}_i; \Theta))^2 \right\}$$

(6)

Θ についての事前分布を $p(\Theta)$ とすると、その事後分布は

$$p(\Theta | \{y_i, \vec{x}_i : i=1, N\}) \propto p(\Theta) \prod_{i=1}^N p(y_i | \vec{x}_i, \Theta)$$

(7)

事後分布を最大にするもの (MAP推定)。

Cu₂O 薄膜結晶の吸収スペクトルの解析

第6回チュートリアルセミナー (2017/11/1) における赤井一郎先生の講演

K. Iwamitsu, S. Aihara, M. Okada and I. Akai, J. Phys. Soc. Jpn. 85, 094716 (2016)

Cu₂O での励起子の Bose-Einstein 凝縮が期待されている。

1. 励起子の寿命が長いことが望まれるので、スピン 3 重項にある、パラ励起子を対象にする。
2. MgO で挟まれた Cu₂O では、僅かな格子不整合により、励起子に対して引力的ポテンシャルが誘起されている可能性が高い。これを定量的に検証する。

吸収スペクトルのベイズ分光

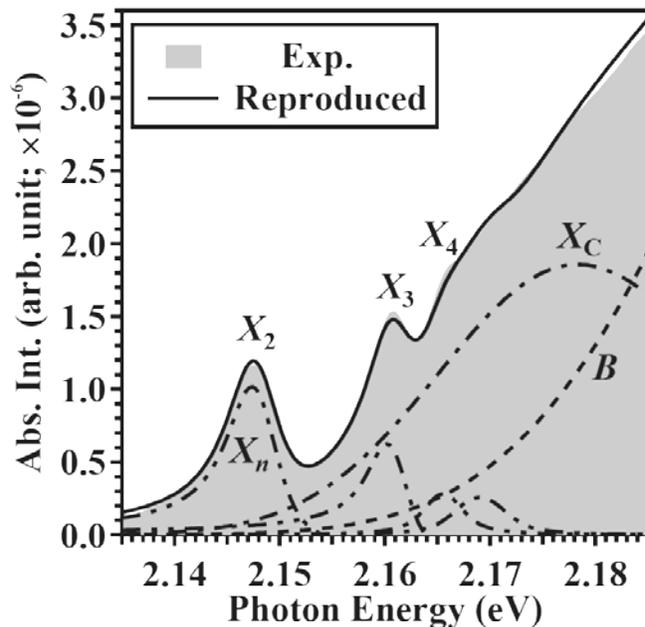
物理モデル: $G(x_i; \Theta)$

➤ X_n 励起子の遷移エネルギー E_n

$$E_n = E_g - \frac{R_y}{n^2}, \quad (n = 2, 3, \dots, n_\infty)$$

➤ X_n 励起子の吸収強度 f_n

$$f_n = f_0 \left(\frac{n^2 - 1}{n^5} \right), \quad (n = 2, 5, \dots, n_\infty)$$



● 励起子ピーク: X_n ($n = 2 \sim 20$)

➤ 非対称ローレンツ形状

✓ 均一幅 Γ_2, Γ_X

$$X_n(E; E_n) = f_n \times \frac{1}{\pi} \frac{\Gamma_n + 2A_n(E - E_n)}{(E - E_n)^2 + \Gamma_n^2}$$

➤ 不均一幅 γ_X

Voigt関数:

$$\int X_n(E; E_n + d\epsilon) \exp\left(-\ln 2 \frac{\epsilon^2}{\gamma_X^2}\right) d\epsilon$$

(X_c, B についても同様)

● バックグラウンド

➤ B : 双極子禁制バンド間吸収

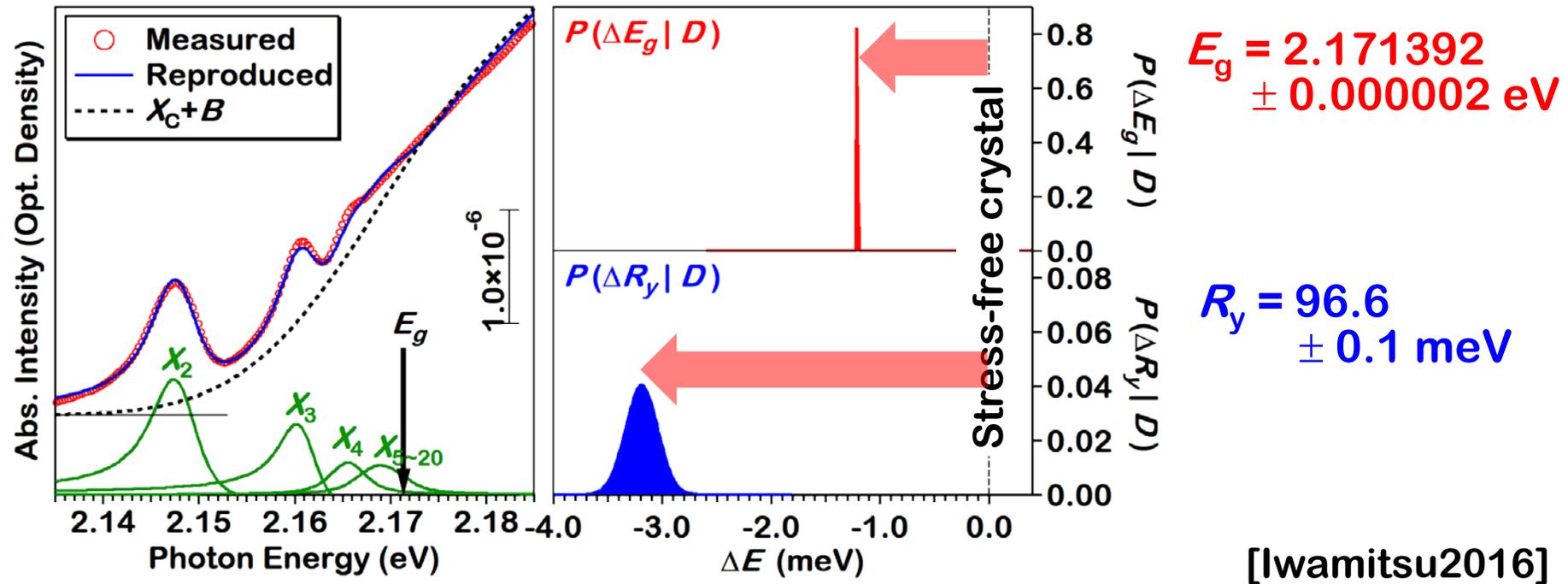
$$X_c(E) = f_{X_c} \pi \beta (1 + \beta^2) \frac{\exp(\pi \beta)}{\sinh(\pi \beta)}, \quad \text{with } \beta = \sqrt{\frac{R_y}{E - E_g}}$$

➤ X_c : 励起子連続帯バンド吸収

$$B(E) = f_B (E - E_g)^{3/2}$$

Burn-inの後、メトロポリス法による250万回のサンプリング

➤ $\sigma_{\text{noise}} = 4 \times 10^{-8}$ (吸収強度) ~ 物理モデル $G(x_i; \Theta)$ のRMSD



- 事後確率分布 $P(\Delta E_g | D)$, $P(\Delta R_y | D)$ は極めてシャープ
- E_g , R_y は明確・有意に変化! → 励起子トラップポテンシャル形成

chemoinformatics:

吉田亮、池端久貴(統数研)、本郷研太(北陸先端大)

H. Ikebata, K. Hongo, T. Isomura, R. Maezono and R. Yoshida,
J. Computer-Aided Mol. Design 31, 379-391 (2017)



定量的構造物性相関 (Quantitative Structure-Property Relation: QSPR)

から

逆定量的構造物性相関 (IQSPR)

へ

分子部品の組み合わせから、自然言語処理の利用への発展

Bayesian inference による

解析・予測から設計へ

の印象的な成果



SPACIER
外挿へ

分子設計については、
津田先生等の MCTS, RNN を用いた ChemTS

DEMO: BAYESIAN MATERIALS DESIGN

- Data: 16,674 organic compounds in PubChem and DFT properties
- Properties: HOMO-LUMO gap, internal energy
- Targeted regions: U_1 , U_2 , U_3

SOFTWARE 'iqspr'

<https://cran.r-project.org/web/packages/iqspr/index.html>

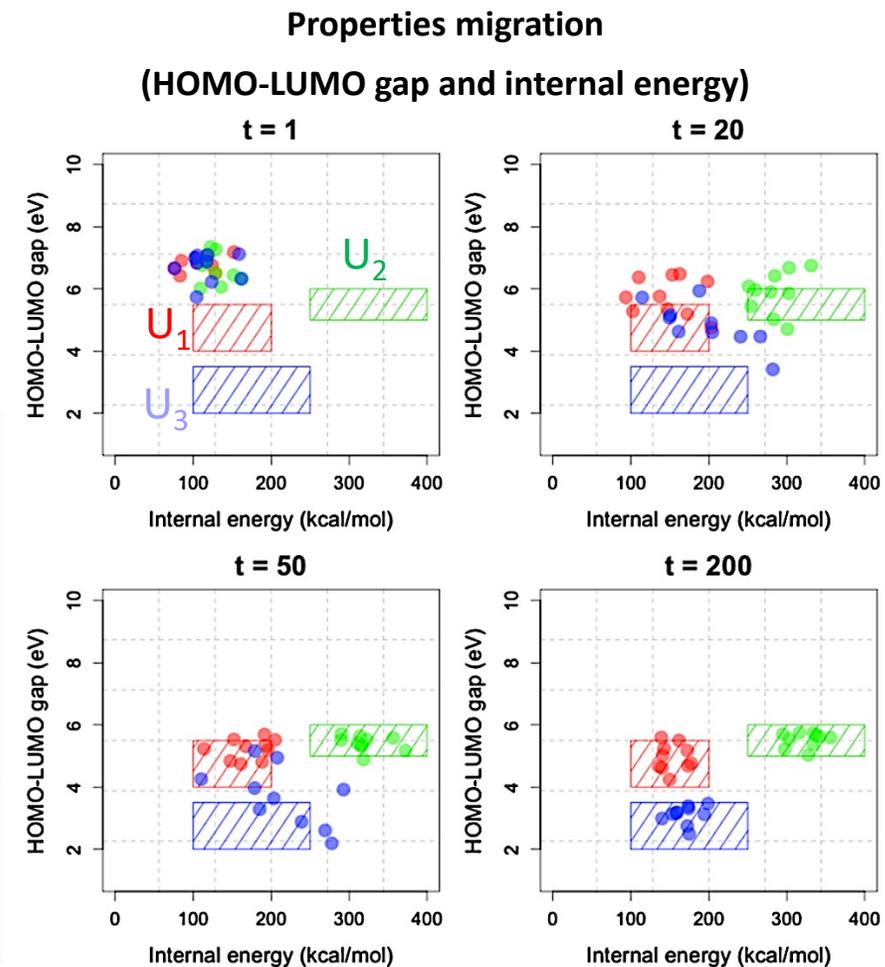
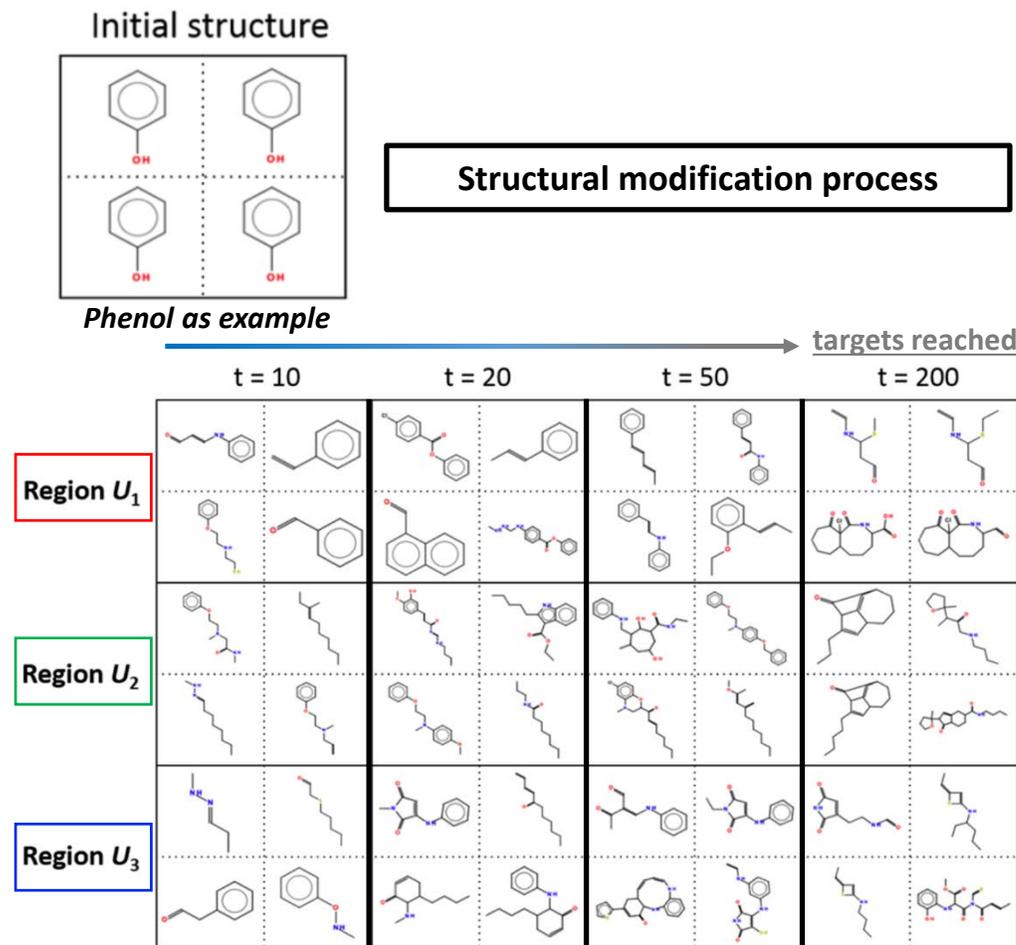


Fig. 4-a (left), 4-b (right) from Ikebata et al., *J. Comput. Aided Mol. Des.*, 2017



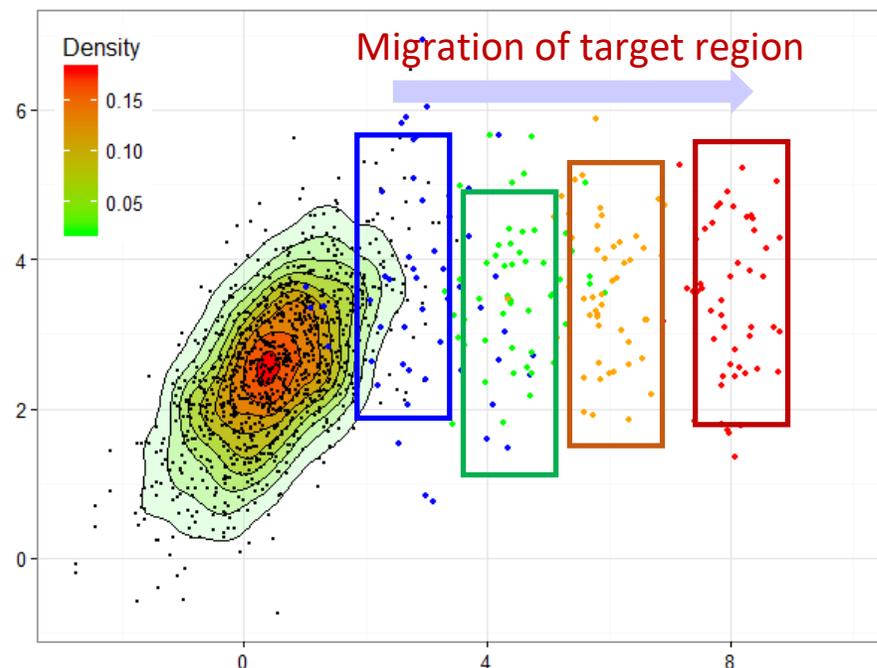
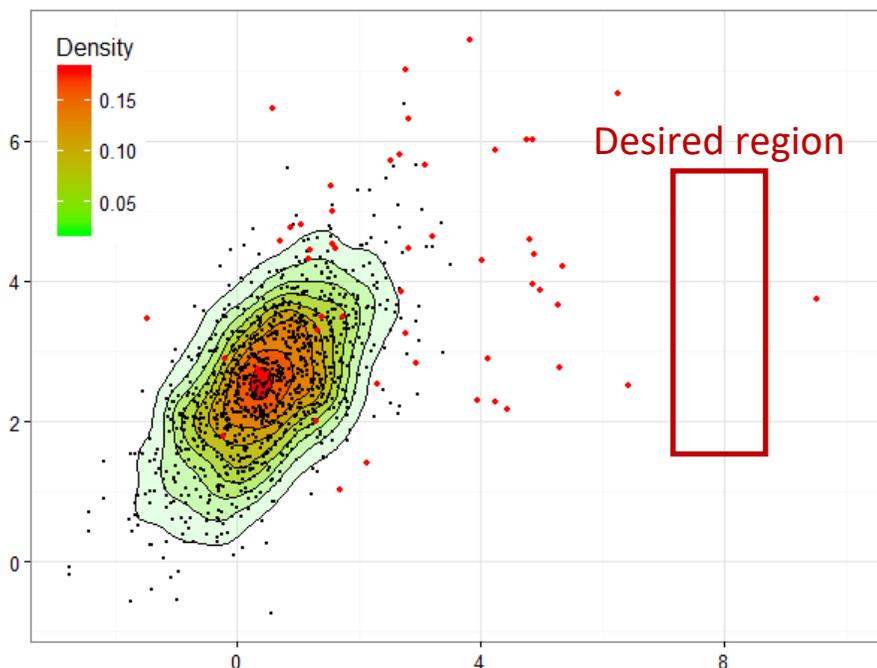
Lambard, G
ISM



SPACIER Go BEYOND INTERPOLATIVE PREDICTION

Difficulty of jumping up to an extreme property region because of exceedingly low accuracy of QSPR models due to the lack of existing materials

Adaptive selection of the target property region by successively controlling the trajectory of created materials



Choice for the temporary target region

Acquisition function of SPACIER

$$A_{U,\lambda}(S, Y^*) = -\|Y^* - \hat{f}(S)\|^2 - \sigma(S) - \lambda \inf_{Y \in U} \|Y^* - Y\|^2$$

Design point for the DFT calculation

転移学習 (Transfer Learning)

関連しているが異なる部分もあるデータから、目的の問題にも利用できる情報・知識だけを取り込んで、より予測精度の高い規則を得る。

神島敏弘 人工知能学会誌 25巻4号 (2010).

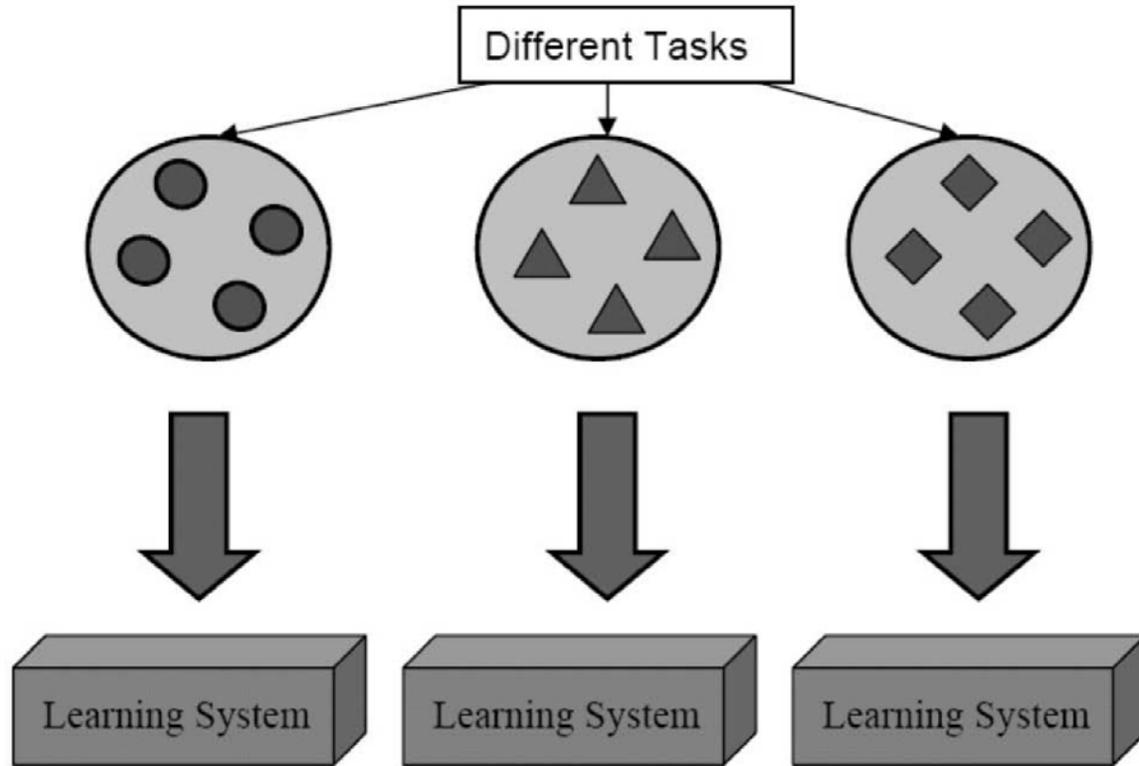
http://www.kamishima.net/archive/2010-s-jsai_tl.pdf

転移学習が役立つ2つのケース

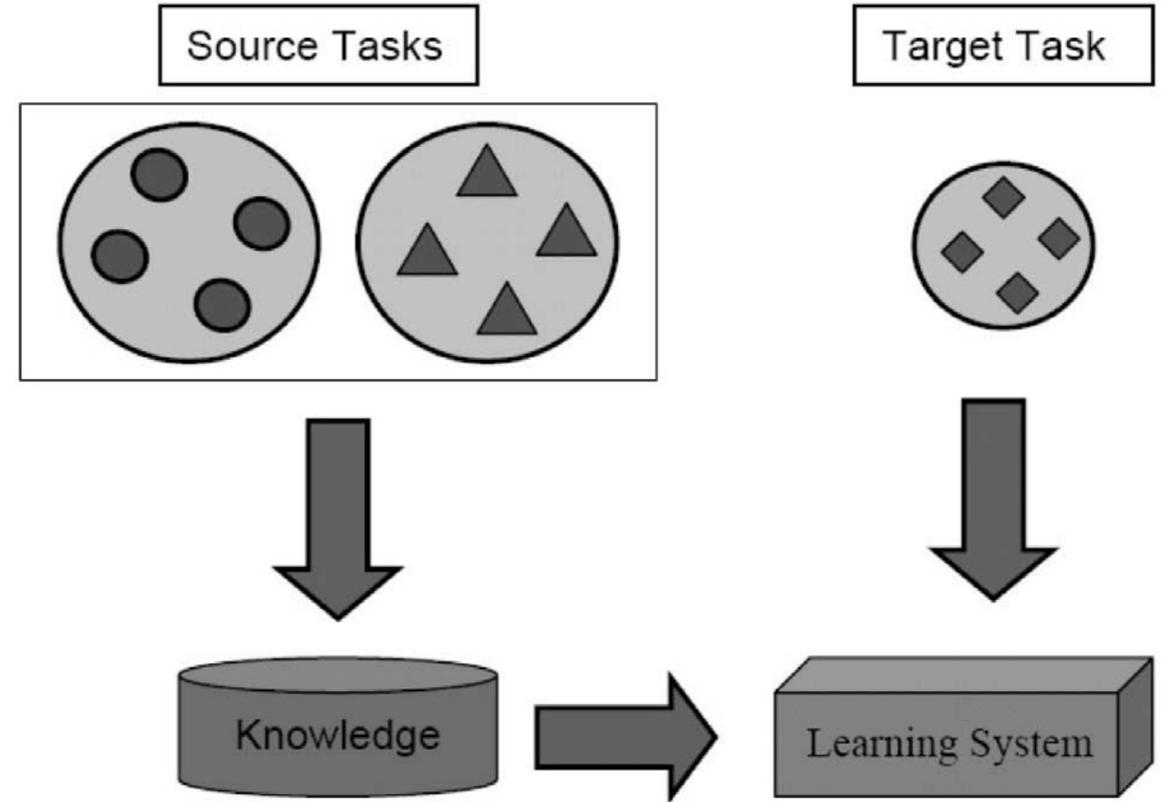
1. 似たような回帰のタスクを行う場合に、回帰の効率と精度を上げる。(GB の例)
むしろ、multi-task と呼ぶほうが適切(配布物に修正)
2. ある目的とする量 A のデータが少ない場合、その量と相関があり、データ取得が容易なために沢山のデータがある量 B というような状況では、量 B の回帰結果を量 A の回帰に使う。(熱伝導の例)

データが得やすい量を通して、データが得にくい量の学習をする。

Learning Process of Traditional Machine Learning



Learning Process of Transfer Learning



S. J. Pan and Q. Yang, A Survey on Transfer Learning,
IEEE Transactions on Knowledge and Data Engineering, 22, 1345 (2010)
https://www.cse.ust.hk/~qyang/Docs/2009/tkde_transfer_learning.pdf

最近、転移学習が物質科学の課題に適用されるようになった。

対応粒界の
安定性の議論

multi-task ?

1. H. Oda, S. Kiyohara, K. Tsuda and T. Mizoguchi,
Transfer Learning to Accelerate **Interface Structure** Searches
J. Phys. Soc. Jpn. 86, 123601 (2017)
2. T. Yonezu, T. Tamura, I. Takeuchi and M. Karasuyama,
Knowledge-Transfer based on Cost-effective Search for **Interface Structures**
arXiv: 1708.03130v1
3. M. L. Hutchinson, E. Antono, B. M. Gibbons, S. Paradiso, J. Ling and B. Meredig,
Overcoming data scarcity with transfer learning
<https://arxiv.org/pdf/1711.05099.pdf>
4. 吉田先生(統数件)らによる、
熱伝導度の予測に、scattering phase space (SPS) のデータを活用
5. 小野先生(KEK)ら
実験の分光スペクトルの解析に、シミュレーションの解析結果を転用

予測と理解

津田先生

(第1回JST workshop (2013/2/11) では)

機械学習の目的は次の2つ

1. 原因究明

2. 予測

それぞれの目的を徹底的に追及する場合は、2つは別の課題。

山本一成氏 (ポナンザの開発者)

『人工知能はどのようにして「名人」を超えたのか?』

第2章: 黒魔術とディープラーニング — 科学からの卒業 —

第1節: 機械学習によってもたらされた「解釈性」と「性能」のトレードオフ

...

最9節: 還元主義的な科学からの卒業

ARTICLE 世界1番の棋士が人工知能に負けた

doi:10.1038/nature.16961

Mastering the game of Go with deep neural networks and tree search

David Silver^{1*}, Aja Huang^{1*}, Chris J. Maddison¹, Arthur Guez¹, Laurent Sifre¹, George van den Driessche¹, Julian Schrittwieser¹, Ioannis Antonoglou¹, Veda Panneershelvam¹, Marc Lanctot¹, Sander Dieleman¹, Dominik Grewe¹, John Nham², Nal Kalchbrenner¹, Ilya Sutskever², Timothy Lillicrap¹, Madeleine Leach¹, Koray Kavukcuoglu¹, Thore Graepel¹ & Demis Hassabis¹

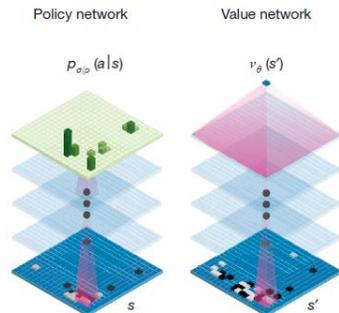
10³⁶⁰

人工知能技術

- **Deep Learning**

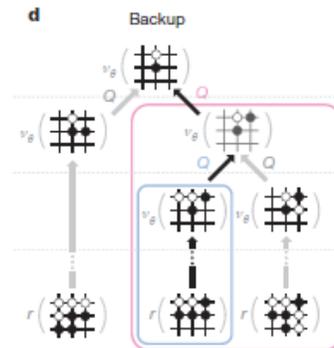
報酬関数 $F(p,v)$

p ; policy network
 v ; value network



機械学習

- **強化学習**
- **モンテカルロ木探索**



棋譜データベース



16万棋譜

3000万局面

計算パワー

Google Cloud
1202CPU
176GPU

人間は絶対に指さない1手目

後手 佐藤天彦 名人

9	8	7	6	5	4	3	2	1
香	桂	銀	金	王	金	銀	桂	香
	飛						角	
歩	歩	歩	歩	歩	歩	歩	歩	歩
歩	歩	歩	歩	歩	歩	歩	歩	歩
	角					金	飛	
香	桂	銀	金	玉		銀	桂	香

先手 PONANZA

1手目 **3八金**

羽生義治

「人間は絶対指さない手ですね。」

中住居に近い形

後手 佐藤天彦 名人

9	8	7	6	5	4	3	2	1
香	桂		金	王		銀	桂	香
	飛	銀				金	角	
歩		歩	歩	歩	歩	歩	歩	歩
	歩						歩	
歩	歩	歩	歩	歩	歩	歩		歩
	角	金		玉		金	飛	
香	桂	銀				銀	桂	香

先手 PONANZA

評価値 -

9手目 5八玉(59)

羽生義治

「この局面になると、・・・『3八金もある手ではないか』と思えるところが、一つの驚きでもあります。」

佐藤天彦名人 対 ポナンザ
第1局 2017年4月

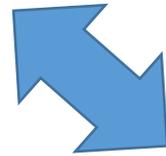
1手目は結局は正しい予測だったが、通常
の理解を超えていた。

中住居に近い

Bayesian inference

少数データへの適応性と解釈性
不確実性への対応

「ひらめき」、「勘」につながる??
STAGE I の課題



Bayesian Deep Learning:
deep learning への不確実性の取り込み
認知 (deep learning) + 推論 (Bayesian inference)
.....

Deep Learning

big data と 解釈不可能性
決定論的

過去の対局の碁譜のデータを入力

碁譜のデータ入力不要

囲碁: アルファ碁、あるいは Master

将棋: Ponanza, あるいは elmo

これらは、人工知能(機械学習)の到達点の例になっている。
強化学習により、データ数が小さいことが問題にならなくなった。

マテリアルズ・インフォマティクスへの水平展開を。

囲碁での探索法(MCTS)のプログラムは津田グループで開発済み