

## 大規模アンサンブル気候データの効率的な解析に向けた コンテンツベース検索システム

中川 友進\*・尾上 洋介\*\*・荒木 文明\*・小山田 耕二\*\*\*

### A Content-Based Database System for Large Ensemble Future Climate Data

Yujin Nakagawa\*, Yosuke Onoue\*\*, Fumiaki Araki\* and Koji Koyamada\*\*\*

**Key words:** Climate Data, Relational Database, Web Application

#### 1. はじめに

コンピュータの性能の向上に伴い、将来気候の大規模なアンサンブルデータ（以後、大規模アンサンブル気候データ）のデータサイズは急激に増加している。例えば、「気候変動リスク情報創生プログラム」により作成された「地球温暖化対策に資するアンサンブル気候予測データベース」(the database for Policy Decision making for Future climate change; d4PDF)<sup>1)</sup>の場合、データサイズは約3 PBである。このような大規模アンサンブル気候データの系統的な解析は、確率的な気候変動の影響評価を行うために非常に重要である。しかし系統的な解析には、一般に高性能なコンピュータだけでなく、大容量のデータ記憶装置が必要である。そのため、個々の研究者による系統的な解析は困難になりつつある。

「気候変動適応技術社会実装プログラム」(Social Implementation Program on Climate Change Adaptation Technology; SI-CAT)は、将来の気候変動に対する適応策に資する技術の開発を目的とした国家プロジェクトである。SI-CATでは、地球の平均気温が産業革命以降2°C上昇した将来気候のアンサンブルデータを作成した<sup>2)</sup>。そのデータはd4PDFの一部として、「データ統合・解析システム」(Data Integration and Analysis System Program; DIAS; <http://www.diasjp.net>)で公開されている。

d4PDFのデータサイズは約3 PBであり、過去のデー

\* 国立研究開発法人海洋研究開発機構  
Japan Agency for Marine-Earth Science and Technology

\*\* 日本大学  
Nihon University

\*\*\* 京都大学  
Kyoto University

タセットと比較して、とても大きい。DIASにおいて、d4PDFが保管されているデータサーバーと同じネットワークに、d4PDFのための解析サーバーがあれば、ユーザーはデータをダウンロードせずに解析を行うことができる。しかし、現状ではd4PDFのための解析サーバーは無いため、ユーザーはd4PDFを自身のローカルコンピュータにダウンロードする必要がある。この場合、d4PDFのデータサイズが大きいことにより以下の問題が発生すると考えている。

- 1) d4PDFをダウンロードするユーザーのディスク容量が不足する。
- 2) DIASのデータサーバーからユーザーのローカルコンピュータへのダウンロードに長い時間かかる。
- 3) データサーバーへ高い負荷がかかる。

これらの問題を解決するためには、ユーザーが要求を満たすデータを検索してダウンロードできる機能を有したユーザーフレンドリーなシステムが必要である。

DIASやWDCC (<https://cera-www.dkrz.de/WDCC/ui/cerasearch>)など、気候データのための従来のWebベースの検索システムは、気候の研究分野において検索や可視化に活用されている。これら全ての従来システムは、データファイルに格納されている物理量を使用してデータを検索するには設計されておらず、データファイルそのものに関連付けられているメタデータを使用してデータファイルを検索するように設計されている。気候データのデータサイズが、ユーザーのローカルコンピュータにダウンロードすることができる程度であれば、これらの従来システムはとても有用である。

本研究では、ユーザーが必要とするデータを検索す

るサービスを提供するために、SI-CATの下で「SI-CAT 気候実験データベースシステム」(System for Efficient content-based retrieval to Analyze Large volume climate data; SEAL)を開発している。SEALを使用すると、物理量など、データファイルの内容に関連付けられているメタデータを使用してユーザーがデータファイルを見つけることができる。そのためSEALは、ユーザーがデータサーバーからローカルコンピュータにダウンロードするデータサイズを大幅に削減できる。

## 2. 設計と実装

### 2.1 データ

SEALの設計に使用したd4PDFは、全球シミュレーションデータ、および領域シミュレーションデータで構成される。全球シミュレーションデータは、気象研究所が開発した水平格子間隔60kmの全球大気モデル(MRI-AGCM<sup>3)</sup>)を用いて作成された。領域シミュレーションデータは、気象研究所によって開発された水平格子間隔20kmの非静力学地域気候モデル(MRI-NHRCM<sup>4,5)</sup>)を用いて作成されており、将来4℃昇温実験、将来2℃昇温実験、過去実験で構成されている。これらのデータのうち、領域シミュレーションデータの地上大気データ<sup>12)</sup>を使用した。地上大気データは「GRIdded Binary」(GRIB)形式で保管されており、時間分解能は1時間である。

### 2.2 基本コンセプト

SEALの設計は、以下の3つの基本コンセプトに基づいて行なった。第1のコンセプトは、ユーザーのニーズを満たすための実用性を持たせることである。ユーザーとして想定されるSI-CATに所属する影響評価の研究者に対してニーズ調査を行なったところ、物理量としては降水量と気温が重要であることが分かった。

また検索条件としては、「ある行政区域の指定期間における、3日間積算降水量の上位から順番に、指定したケース数だけ表示する。」などが重要であることが分かった。第2のコンセプトは、生データ(GRIB形式のデータ)に格納されている物理量を出来るだけ加工しないことである。そのため、バイアス補正など、ユーザーに依存する加工は行わないこととした。第3のコンセプトは、SEALのために開発された技術を様々なデータセットに適用するための汎用性である。

図1にSEALの概念図を示す。SEALは、リレーショナルデータベース、データ提供機能、およびユーザーインターフェースで構成されている。これら3つ機能のうち、時間的、空間的に圧縮されたデータを登録するように設計されているPostgreSQLを使用したリレーショナルデータベースが重要な役割を担っている。データ提供機能は、ユーザーが検索結果に基づいて時間的、空間的に切り出されたデータをダウンロードする機能を提供する。さらに、Webベースのユーザーインターフェースにより、ユーザーはPostgreSQLの知識がなくてもリレーショナルデータベースを簡単に使用できる。

### 2.3 リレーショナルデータベースの設計

ユーザーは、自身が研究の対象としている特定の領域に関心があるため、単一の格子点ではなく、複数の格子点の組み合わせによって定義される領域(行政区域や流域など)における物理量を必要としている。さらに、ユーザーは、気温のような物理量については、時別値ではなく、日別値、月別値または年別値を必要とする。そこで、d4PDFの物理量に対して、空間的、時間的な圧縮を行い、それらの圧縮された物理量をリレーショナルデータベースへ登録することとした。ここで、空間的圧縮とは、複数の格子点における物理量

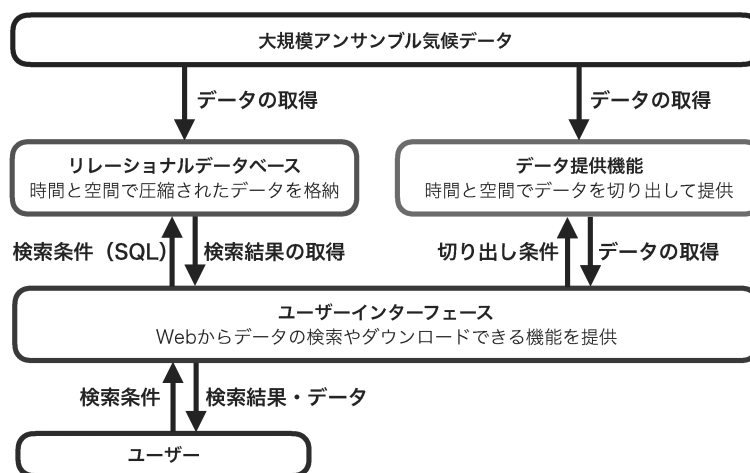


図1 SEALの概念図。

を平均することである。また、時間的圧縮とは、特別値を積算（例えば降水量の場合）または平均（例えば気温の場合）することにより、日別値、月別値、年別値を計算することである。

## 2.4 データ提供機能の設計

データ提供機能により、ユーザーは、リレーショナルデータベースによる検索結果に基づいて、時間的、空間的に切り出された生データをダウンロードすることができる。DIASは、ユーザーが「Grid Analysis and Display System」(GrADS)のバイナリデータ（ビッグエンディアン、ヘッダーなしの4バイト浮動小数点）をダウンロードできる機能を提供している。しかし、ほとんどのユーザーは、テキスト形式やCSV形式などの人間が読める形式へ変換したデータを必要としている。そこで、GrADSバイナリデータをテキスト形式またはCSV形式へ変換する機能をPythonスクリプトで開発した。

## 2.5 ユーザーインターフェースの設計

ユーザーからのニーズを考えると、ユーザーがGRIB形式やPostgreSQLに関する知識がなくても検索を行えるシステムであることが望ましい。また、コマンドラインインターフェースの知識がなくてもGrADSバイナリデータをテキスト形式またはCSV形式へ変換できることが望ましい。そこで我々は、これらの知識がなくてもリレーショナルデータベースとデータ提供機能を利用できるWebベースのユーザーインターフェース（Web UI）を開発した。Web UIは、条件に合うデータを見付けることを主な機能とした「SEAL-Finder」(SEAL-F)、データを分析および可視化することを主な機能とした「SEAL-Visualizer」(SEAL-V)で構成されている。そのため、ユーザーは、SEAL-Fを使用すると検索の返り値を数値として取得することができ、またSEAL-Vを使用すると検索の返り値をヒストグラムなどの画像として取得することができる。

## 2.6 リレーショナルデータベースの実装

検索の高速化と高い実用性の両方を満たすために、降水量は、特別値と日別値の2つの異なる時間分解能のデータをリレーショナルデータベースへ登録することとした。また、気温は、日平均値、日最大値、日最小値を登録することとした。ほとんどの場合、ユーザーは自身が研究の対象とする行政区域の物理量を必要とする。本研究で使用している地上大気データの空間分解能が20 kmであることを考慮して、47都道府県それぞれについて、行政区域と重なる格子点の物理量を積算して、リレーショナルデータベースへ登録するこ

ととした。

## 2.7 ケーススタディ

### 2.7.1 空間的圧縮、時間的圧縮による圧縮率

空間的圧縮、時間的圧縮を適用した際に、生データと比べて、どのくらいデータサイズが縮小されるかを定量的に調べるため、それらの圧縮率を計算した。空間的圧縮の圧縮率は、北海道（302個の格子点で構成）では $1/302 \approx 0.33\%$ 、東京都と大阪府（それぞれ13個の格子点で構成）では $1/13 \approx 7.7\%$ となる。また、時間的圧縮の圧縮率は、特別値から日別値へ圧縮した場合は $1/24 \approx 4.2\%$ となる。そのため、空間的圧縮と時間的圧縮の両方を適用した場合、日別値のデータサイズは、最大で0.01%、最小で0.3%へ縮小する。データサイズの縮小は、従来システムを使用する方法と比較して、必要なデータファイルを取得するまでの所要時間の短縮に繋がる。

### 2.7.2 必要なデータファイルの取得までの所要時間

SEALの利点を明確にするために、従来システムを使用した方法と国立研究開発法人海洋研究開発機構（Japan Agency for Marine-Earth Science and Technology; JAMSTEC）のローカルサーバー上のSEALを使用した方法について、必要なデータファイルを取得するまでの所要時間を定量的に調べた。ローカルサーバーは、Intel Xeon E7-4820（CPU; 40コア）と512 GBの物理メモリを搭載している。すべての解析において1つのCPUコアを使用した。図2はSEALの利点の概略図であり、北海道で降水量の日別値が100 mmを超える日について、将来4°C昇温実験の降水量と気温の特別値のデータファイルをダウンロードする状況を想定している。従来システムを使用した方法では、北海道周辺の降水量の特別値のデータファイルのダウンロードに3~90時間、条件を満たす日にちの探索に約31時間、条件を満たす日にちにおける気温の特別値のデータファイルのダウンロードに1~30分を要する。一方、SEALを使用した方法では、条件を満たす日にちを約5秒で見付けることができ、条件を満たす日にちの降水量と気温の特別値のデータファイルを2~60分でダウンロードできる。ダウンロードする特別値のデータファイルのデータサイズは、生データと比較して約0.5%へ減少する。したがって、従来のシステムを使用した方法と比較して、SEALは必要なデータファイルを取得するまでの所要時間を1%未満へ減らすことができる。

### 2.7.3 大雨の検索

SEALの実用性を明確にするために、河川工学を専門とするユーザーが利用することが多いと考えられる

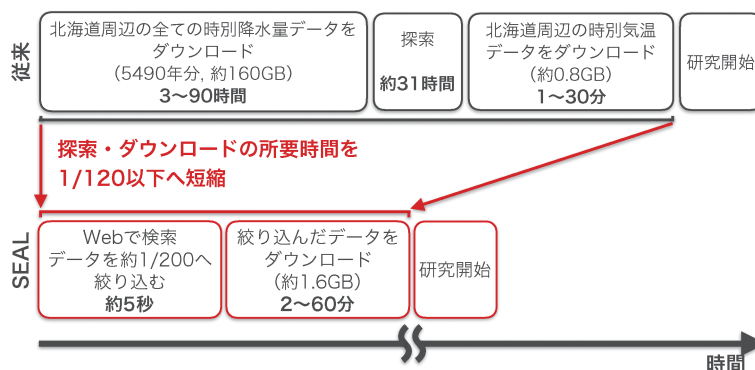


図2 SEALの予想される利点の概略図。図に示されている経過時間は、DIASからJAMSTECへのダウンロード速度(4~120 Mbps)に基づく、将来4度昇温実験のデータを使用した際の実測値である。

検索条件を想定し、ローカルサーバー上で動作するSEALを用いて検索を行った。将来4℃昇温実験の東京都のデータのうち、連続降水日数が20日を超え、かつ累積降水量が800 mmを超えるという検索条件を満たすイベントを抽出した。ここで、降水量の日別値が0.1 mmを超える日を降水日と定義している。その結果、検索に約30秒かかり、条件を満たす6つのイベントが見つかった。その中でも、累積降水量が最大のイベントは、24日間に約1,109 mmを示した(図3)。累積降水量の等高線図を調べたところ、このイベントは、静岡県を中心とする大雨であった。従来の方法で上記の現象を調べた場合、東京都に対応する降水量の特別値のデータファイルのダウンロードに0.13~3.9時間かかる。さらに、ユーザーは自身で連続降水日数と累積降水量を計算する必要があるため、さらに時間がかかる。したがって、SEALは必要なデータファイルを取得するまでの所要時間を短縮する十分な能力を持っている。

### 3. 考察

従来のフレームワークの1つであるOPeNDAP(<https://www.opendap.org>)は、単一の時系列グリッドデータ(すなわち、単一のデータファイル)から値を検索する機能を提供している。しかし、データファイル内のすべての値をスキャンするため、検索速度は速くはない。一方、SEALは、リレーショナルデータベースを採用することにより、複数のアンサンブル(すなわち全てのデータファイル)から基準を満たす値を高速に検索できる。これは、データベースインデックスの影響を強く受ける検索については、リレーショナルデータベースは検索の速度が大幅に向上するためである。さらにリレーショナルデータベースは、複数のアンサンブルを一度にスキャンできるという利点を有

する。

第1章で述べたように、大規模アンサンブル気候データから必要なデータファイルを取得するには3つの問題が予想される。すなわち、ユーザーのディスク容量の不足、長時間のダウンロード、データサーバーへの高い負荷、である。第2.7節のケーススタディにおいて、ユーザーがローカルコンピュータにデータファイルをダウンロードするデータサイズおよび所要時間は、従来の方法と比較して、それぞれ約0.5%および約1%に縮小された。その結果、1番目と2番目の問題が解決された。さらに、データサイズと必要なデータを取得するまでの所要時間が大幅に削減されるため、データサーバーへの負荷が大幅に減少した。したがって、3番目の問題も解決された。その結果、SEALは予想どおりに機能し、大規模アンサンブル気候データから必要なデータを取得する際に予想される問題をすべて解決した。

### 4. 教訓

第2.2節で述べたように、SEALの開発において、ユーザーのニーズを満たす実用性が重要なコンセプトの一つである。そこで、SEALの開発を始めるにあたり、ユーザーとして想定される研究者へのニーズ調査を行なった。ニーズ調査に際して、回答の自由度が高すぎると、SEALでは実現できない要求が出てくる可能性があり、一方で制限を付けすぎると回答を誘導する可能性がある。そのため、正しくニーズを引き出すためには、アンケートへどのように設問を記載するかという点が重要となる。SEALの検索において、検索条件、および「どのような戻り値が必要となるか?」という検索結果を問う設問は重要であり、特に慎重に議論を行った。その結果、こちらが想定している具体的な回答例を併記することにより、知りたいニーズを

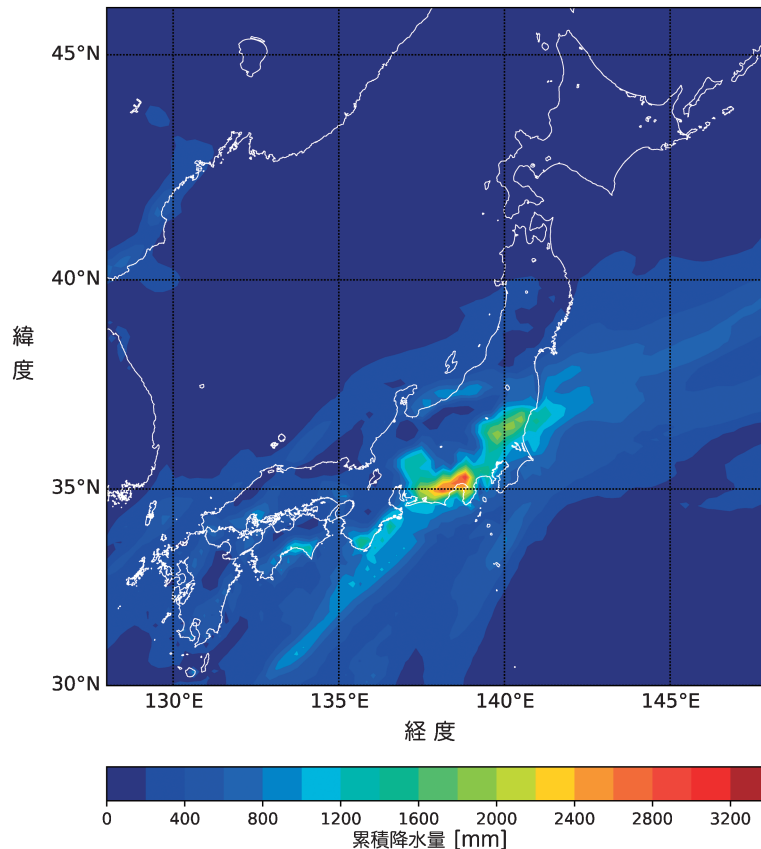


図3 SEALを使用して見付けた大雨の累積降水量の等高線図。

得ることができた。ただし戻り値については、d4PDFに格納されていないデータへの要望があったため、戻り値として可能な範囲を示すべきだったと考えている。すなわち、d4PDFの生データに対して基本的に加工は行わないこと（第2.2節を参照）を明確しておくべきだったと考えている。SEALの開発のチームメンバーだけでなく、他のSI-CATのメンバーの意見を聞くことにより、設問の問題点を洗い出すことができ、ユーザーのニーズを正しく引き出し、ニーズを満たすシステムを開発できたと考えている。

## 5. 結論

大規模アンサンブル気候データのデータサイズが増加するにつれて、従来のWebベースの検索システムを使用した検索に際して、3つの問題が予想される。すなわち、ユーザーのディスク容量の不足、長時間のダウンロード、データサーバーへの高い負荷、である。これらの問題を解決するために、物理量などのデータファイルの内容に関連付けられたメタデータを使用してユーザーがデータファイルを検索できるようにするシステムであるSEALを開発した。SEALの使用に際して、ユーザーはGRIB形式とPostgreSQLの知識は

必要とせず、ユーザーは検索結果に基づいてGrADSバイナリ形式、テキスト形式、CSV形式のデータをダウンロードすることができる。さらに、ユーザーはWebベースのユーザーインターフェースを介してすべての作業を行うことができる。SEALのリレーショナルデータベースに格納されるデータサイズは、空間的圧縮および時間的圧縮を採用することにより、生データと比較して最大で0.01%、最小で0.3%に縮小される。リレーショナルデータベースは、データベースインデックスの影響を大きく受ける検索に対しては、検索速度を大幅に改善すると共に、複数のアンサンブルを一度にスキャンすることができる。さらに必要なデータの取得に際して、データサイズと所要時間をそれぞれ約0.5%と約1%へ減らすことができる。これらの利点により、データサーバーへの負荷が削減される。その結果、SEALは期待通りに動作し、すべての問題を解決することができた。

SEALは現在、JAMSTECのローカルサーバーにおいて運用試験中であり、2019年度末までにDIASから公開される予定である。SEALのために開発された技術は、他の研究分野におけるシミュレーションや観測のデータに対しても、非常に有用だと考えている。

## 謝辞

本研究は、気候変動適応技術社会実装プログラム (SI-CAT) の一環として行われている。また、本研究は SOUSEI と SI-CAT の下で作成された、地球温暖化対策に資するアンサンブル気候予測データベース (d4PDF) を使用した。

## 参 考 文 献

- 1) R. Mizuta, A. Murata, M. Ishii, H. Shiogama, K. Hibino, N. Mori, O. Arakawa, Y. Imada, K. Yoshida, T. Aoyagi, H. Kawase, M. Mori, Y. Okada, T. Shimura, T. Nagatomo, M. Ikeda, H. Endo, M. Nosaka, M. Arai, C. Takahashi, K. Tanaka, T. Takemi, Y. Tachikawa, K. Temur, Y. Kamae, M. Watanabe, H. Sasaki, A. Kitoh, I. Takayabu, E. Nakakita, and M. Kimoto: Over 5,000 years of ensemble future climate simulations by 60-km global and 20-km regional atmospheric models, *Bulletin of the American Meteorological Society*, vol. 98, no. 7, pp. 1383–1398 (2017).
- 2) M. Fujita, R. Mizuta, M. Ishii, H. Endo, T. Sato, Y. Okada, S. Kawazoe, S. Sugimoto, K. Ishihara, and S. Watanabe: Precipitation changes in a climate with 2-k surface warming from large ensemble simulations using 60-km global and 20-km regional atmospheric models, *Geophysical Research Letters*, vol. 46, no. 1, pp. 435–442 (2019).
- 3) R. Mizuta, H. Yoshimura, H. Murakami, M. Matsueda, H. Endo, T. Ose, K. Kamiguchi, M. Hosaka, M. Sugi, S. Yukimoto, S. Kusunoki, and A. Kitoh: Climate simulations using mri-agcm3.2 with 20-km grid, *Journal of the Meteorological Society of Japan. Ser. II*, vol. 90A, pp. 233–258 (2012).
- 4) H. Sasaki, A. Murata, M. Hanafusa, M. Oh'izumi, and K. Kurihara: Reproducibility of present climate in a non-hydrostatic regional climate model nested within an atmosphere general circulation model, *SOLA*, vol. 7, pp. 173–176 (2011).
- 5) A. Murata, H. Sasaki, M. Hanafusa, and K. Kurihara: Estimation of urban heat island intensity using biases in surface air temperature simulated by a nonhydrostatic regional climate model, *Theoretical and Applied Climatology*, vol. 112, pp. 351–361 (2013).