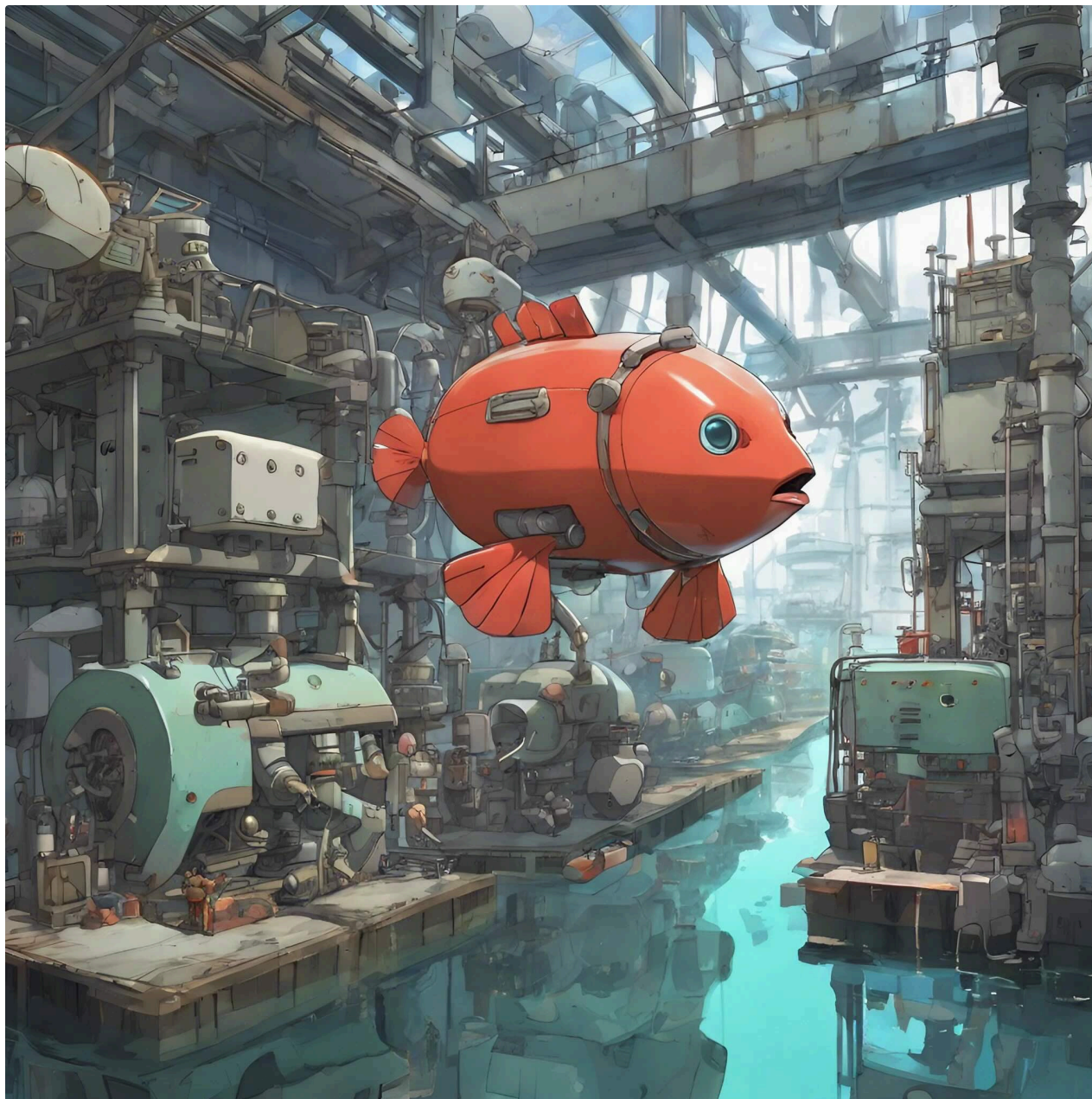


sakana.ai

世界初、100%AI生成の論文が査読通過 「AIサイエンティスト」が達成

March 13, 2025



注：この実験はICLR主催者とそのワークショップ関係者の全面的な協力のもとで実施されたものです。「透明性と倫理的行動規範の重要性」のセクションを参照。

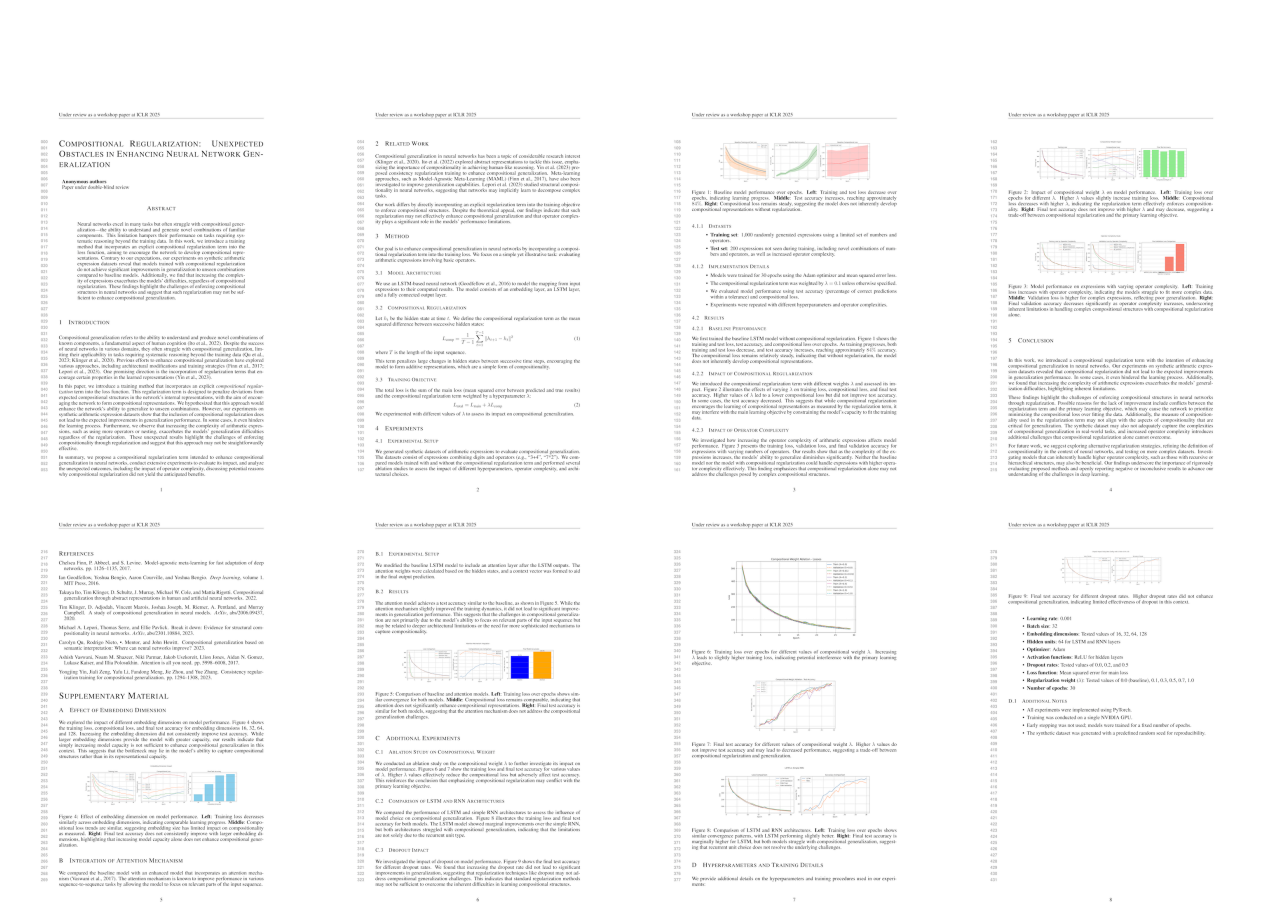
「[AIサイエンティスト](#)」によって作成された論文が、トップレベルの機械学習学会のワークショップで査読プロセスを通過しました。私たちの知る限り、これは完全AI生成論文が査読プロセスを通過した初めての事例です。

この論文は、昨年発表した[AI Scientist](#)の改良版であるAI Scientist-v2によって生成されたものです（AI Scientist-v2の詳細は今後公開予定）。私たちは、国際学会ICLR 2025で行われる[ワークショップ](#)の協力を得て、AI生成論文を二重盲検レビューのプロセスに提出する実験を行いました。このワークショップは現実世界における深層学習の限界と課題に焦点を当てたものであり、幅広い研究を受

け付けていることから今回の実験に適したものとして選ばれました。ICLRは、NeurIPSおよびICMLと並ぶ、機械学習とAI研究における世界の三大国際学会の1つです。

この実験は、ICLR主催者とそのワークショップ関係者の全面的な協力のもとで実施されました。また、ブリティッシュコロンビア大学からは、本研究に関するIRB（研究倫理審査委員会）の承認も得ました。この実験に協力してくださった皆様に感謝申し上げます。本研究についてはICLRワークショップで発表を行い、本プロジェクトで私たちが得た経験や課題を共有する予定です。

本研究はブリティッシュコロンビア大学およびオックスフォード大学のチームとの協力のもと行われました。



トップレベルの国際的なAI学会のワークショップで査読プロセスを通過した100%AI生成による論文。

評価プロセス

ICLRワークショップの主催者との調整のもと、私たちは3つのAI生成論文をワークショップに提出しました。査読者の負担を考慮し3つの論文を査読にかけることをワークショップの主催者と合意しました。査読者には、査読している論文がAIによって生成されたものでありうること（43件中3件）が知らされていましたが、どの論文が実際にAI生成であるかは知らされませんでした（査読プロセスの詳細については、ICLRワークショップの査読者ガイドライン参照）。

これらのAI生成論文は、AIによって完全にエンドツーエンドで生成され、人間の修正を加えていません。AI Scientist-v2は、科学的な仮説立案から、仮説検証のための実験の考案、実験を実施するためのコードの作成と改良、実験の実行、データ分析、データの視覚化、タイトル・参考文献・図の配置・書式設定を含む論文原稿執筆のすべてのプロセスを自動で行いました。人が行うのは、ワークショップのテーマに合うよう大まかなトピックを与えるのみです。

AIが生成した論文から、内容の多様性と質を両方考慮しつつ3つの論文を選び提出しました（3つの論文については、後述のように詳細な内部レビューを行いました）。その結果、2つの論文は採択基準を至らず、残りの1つは査読者の平均スコア6.33で採択水準を上回りました。これはワークショップに提出された全論文の約45%の水準であり、多くの研究者の論文を上回るものです。査読者による具体的なスコアは次のとおりです。

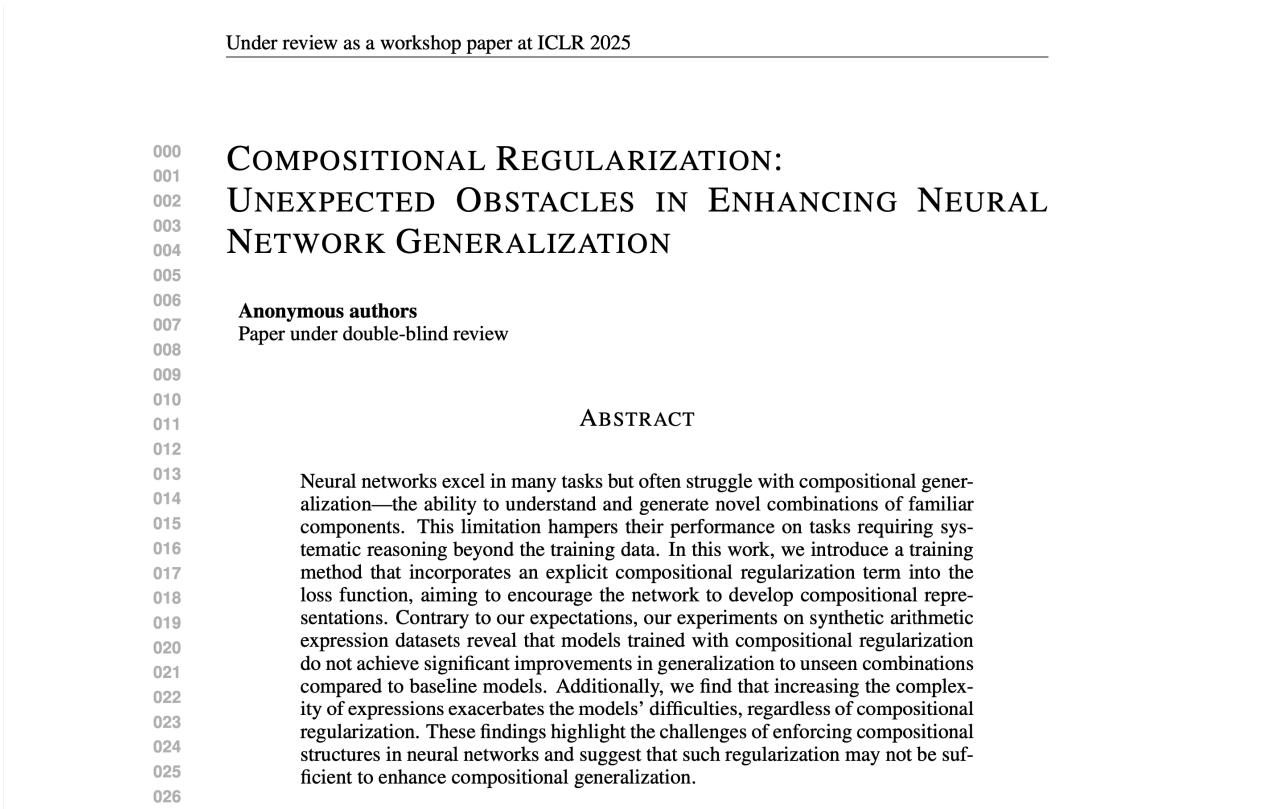
- 6点: 採択水準ぎりぎりで合格
- 7点: 良い論文、採択

- 6点：採択水準ぎりぎりで合格

ただし、今回の実験では、AI Scientistによる論文が受け入れられたとしても、実際に公開される前に撤回されることが事前に合意されていました。これは、AI生成論文を、人間の論文と同じ場で扱って良いかに関する研究者コミュニティのコンセンサスはまだ存在していないためです。

なお、本論文は査読プロセスの後に撤回されたため、ICLRワークショップの主催者による追加のメタレビューは実施されませんでした。査読者による平均スコア6.33にもかかわらず、メタレビューアがこの論文を不採択とした可能性は理論上はあることを付記しておきます。

昨年発表したAI Scientistの初版は、AIが科学論文を丸ごと生成できることを初めて示したものでした。今回の結果（AI Scientist 「第2版 (v2)」）は、100%AIで生成された論文が標準的な科学的査読プロセスを通過するに足る水準に達した、私たちが知る限り世界で初めての成果です。



AI Scientist-v2が、与えられた大まかな研究トピックだけをもとに生成した論文「Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization」。内容は、ニューラルネットワークの構成的汎化能力を高めるために新たな正則化手法を検討した際に遭遇した障害について報告したもの。この論文はICLRワークショップにて採択水準を上回る平均査読者スコア6.33を得た。

透明性と倫理的行動規範の重要性

科学コミュニティにとって、AI生成研究の質がどこまでに達しているのかを知ることは極めて重要であり、その最良の方法の1つが、人間のための評価手法である厳格な査読プロセスにサンプルを提出していただくことです（ただし、そのプロセスを管理する人々の許可を得ていることが条件です）。

前述のように、この研究はICLR主催者とそのワークショップ関係者の全面的な協力のもとで実施されました。AI生成論文の査読プロセスでの評価に関する本研究にご協力いただいた皆様に感謝します。また、本研究においてはブリティッシュコロンビア大学からIRBの承認も得ています。

なお、今回のAI生成論文をICLRワークショップのOpenReviewフォーラムで公開する予定はありません。これは、今回の実験の目的に照らし、ICLR主催者とワークショップ関係者との合意で、AI生成論文は査読プロセスが完了した後に自動的にデスクリジェクトされるとしていたためです。

私たち科学コミュニティとして、AIが生成する科学に関する規範を形作っていくことが重要です。たとえば、AIが全体的あるいは部分的に生成した論文について、いつ、どのようにそれを開示すべ

069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085

3 METHOD

Our goal is to enhance compositional generalization in neural networks by incorporating a compositional regularization term into the training loss. We focus on a simple yet illustrative task: evaluating arithmetic expressions involving basic operators.

3.1 MODEL ARCHITECTURE

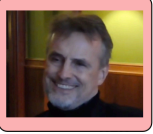
We use an LSTM-based neural network (Goodfellow et al., 2016) to model the mapping from input expressions to their computed results. The model consists of an embedding layer, an LSTM layer, and a fully connected output layer.

3.2 COMPOSITIONAL REGULARIZATION

Let h_t be the hidden state at time t . We define the compositional regularization term as the mean squared difference between successive hidden states:

$$\frac{1}{T-1} \sum_{t=2}^T \|h_t - h_{t-1}\|^2$$

Comment:
This should be Hochreiter & Schmidhuber



Comment:
This should be more precise.

AI Scientistは時々、引用の間違いを犯した。ここでは「LSTMベースのニューラルネットワーク」の出典をGoodfellow (2016) としていたが、正しくはHochreiter & Schmidhuber (1997)。

レビューとコメントに加え、レビューの初期段階では各論文にNeurIPSやICLRのようなトップML学会のガイドラインに沿ったスコアをつけました。

さらに、AI Scientist-v2によって作成された実験結果が再現可能であることを確認するために、コードレビューも実施し、図の欠落、引用の過度の欠落、書式設定の問題などのエラーをチェックしました。なお、結果の科学的な正確性、再現性、統計的な厳密性を向上させるために、AI Scientistは論文に含める各実験を数回繰り返す努力をするようデザインされています。

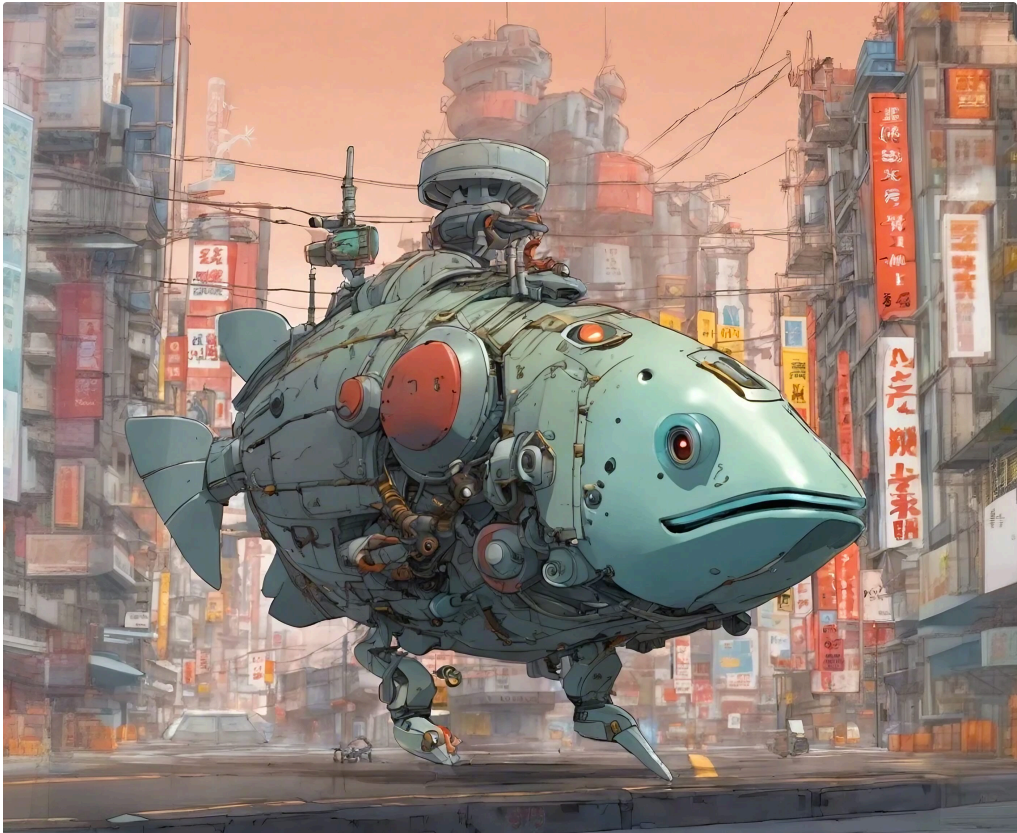
結果的に、3つの論文のいずれも、そのままの形では、ICLRの本会議の論文としては私たちが考える基準を満たしていないと結論付けました。とはいえ、ワークショップに提出した3つの論文には、それぞれ興味深く、独創的なアイデアが含まれており、暫定的ながらも発展する余地があるものだと考えました。したがって、ICLRワークショップには適格である可能性があると判断しました。

GitHubリポジトリにて、これらの3つのAI生成論文とともに、私たち自身の人間によるレビューを公開しています。ぜひこれらの論文を実際に見ていただき、ご自身でも評価していただければと思います。

AIサイエンティストの未来

次の世代のAIサイエンティストは、新しい科学の時代を切り開くと確信しています。AIがトップレベルの機械学習の国際学会ワークショップで査読を通過する論文を丸ごと生成できるという事実は、これからの進歩の確かな兆候です。しかし、これはほんの始まりに過ぎません。AIは今後も改善し続け、もしかしたらそれは指数関数的な向上かもしれません。未来のある時点で、AIはおそらく人間のレベルと同等かそれ以上の論文を生成できるようになるでしょう。AIサイエンティストのようなシステムは、機械学習のトップ学会だけでなく、科学のトップジャーナル（学術誌）でも受け入れられる水準の論文を生成するようになると私たちは予想しています。

もちろん、AI科学と人間科学を比べることが最終目的ではありません。最も大事なものは、人間やAIによる科学がもたらす発見が、病気の治療につながったり、宇宙を支配する法則を明らかにしたりするなど、人間の繁栄に役立つことでしょう。AI科学には、人間社会をよりよく変化させる可能性があります。その時代を切り開く一端を担えることを、楽しみにしています。



Sakana AI

日本でのAIの未来を、Sakana AIと一緒に切り拓いてくださる方を募集しています。当社の[募集要項](#)をご覧ください。

© Sakana AI 株式会社

