

# MLP の重みを反映した Sparse Autoencoder の初期化手法の提案

菊谷 幹<sup>1</sup> 北田 俊輔<sup>1</sup> 原 聡<sup>1</sup><sup>1</sup> 電気通信大学

{k2210202@gl.cc., ka004763@gl.cc., satohara@}uec.ac.jp

## 概要

Sparse Autoencoder (SAE) は LLM 内部の処理を解釈するためのツールとして広く用いられている。しかし、異なる初期値で訓練した SAE では異なる“特徴”が学習されることが報告されている。この問題を解決して SAE の学習を安定化させるためには、SAE の良い初期値を設定する必要がある。本研究では SAE の初期値として Transformer 層の MLP のデコーダを用いる手法として SAE-MD を提案する。実験の結果、SAE-MD は既存の SAE と同等の性能を有し、学習の安定性を向上できることを確認した。

## 1 はじめに

大規模言語モデル (LLM) の発展に伴い、LLM の出力に対する安全性の追求やさらなる改良のため、モデル内部での情報の処理方法を明らかにする試みが数多くなされている [1, 2]。そのような試みの一つが Sparse Autoencoder (SAE) である [3, 4, 5]。SAE は、LLM 内部の情報を“特徴” (≈ 人間の解釈可能な概念) としてニューロン一つ一つに対応付けることを目的とした手法である。SAE のニューロンの発火の様子を見ることで、SAE が獲得した“特徴”を用いて LLM の内部の処理を解釈できるようになる。

しかし、異なる初期値で訓練した SAE では異なる“特徴”が学習されることが報告されている [6]。この事実は SAE がいくつかの“特徴”を捉えきれない、あるいは逆にノイズなどの無関係な“特徴”を捉えてしまっている可能性を示唆している。このような不安定な SAE を用いて LLM の挙動を解釈しても、その解釈結果が不正確である懸念が残る。SAE が LLM の主要な解釈手段の一つとなっている現在において、その不安定性の解消は喫緊の課題である。

SAE の初期値への強い依存性を解消し、“特徴”の学習を安定化させるためには、SAE の良い初期値を

設定する必要がある。そこで、本研究では SAE の初期値として Transformer 層の MLP のデコーダを用いることを提案する。本研究の着眼点は、MLP のデコーダのベクトルの和として記述されている LLM の埋め込み表現を、SAE は改めて複数のベクトルの和として記述し直していると解釈する点にある。この事実は MLP のデコーダが SAE のデコーダとして有用であることを示唆している。実験の結果、MLP のデコーダを初期値に用いた SAE は既存の SAE と同等の性能を有しつつも、学習の安定性を大きく向上できることを確認した。

## 2 Sparse Autoencoder (SAE)

SAE は入力ベクトルをスパースなベクトルに変換し、そのベクトルを使って入力ベクトルを再構成するように学習するモデルである [3]。SAE に LLM 内部の埋め込み表現を入力することでスパースなベクトルに変換できる。このとき、スパースなベクトルの各要素が“特徴”に対応する。ベクトルの  $i$  番目の要素が発火し非ゼロであれば、その要素に対応する“特徴”が埋め込み表現中に存在すると解釈できる。

様々な構造の SAE が提案されており、代表的な SAE の一つが活性化関数に JumpReLU を用いた JumpReLU SAE である [7]。 $\mathbf{x} \in \mathbb{R}^d$  を SAE への入力ベクトルとしたとき、 $n$  次元の中間層を持つ JumpReLU SAE は式 (1), (2), (3) で表される。

$$\mathbf{z} = \text{JumpReLU}_\theta(\mathbf{W}^{\text{enc}}(\mathbf{x} - \mathbf{b})), \quad (1)$$

$$\hat{\mathbf{x}} = \mathbf{W}^{\text{dec}}\mathbf{z} + \mathbf{b} = \sum_i z_i \mathbf{w}_i^{\text{dec}} + \mathbf{b}, \quad (2)$$

$$\text{JumpReLU}_\theta(z) = z \text{ (if } z > \theta), \quad 0 \text{ (otherwise)}. \quad (3)$$

$\mathbf{W}^{\text{enc}} \in \mathbb{R}^{n \times d}$  と  $\mathbf{W}^{\text{dec}} \in \mathbb{R}^{d \times n}$  は SAE のエンコーダとデコーダの重み行列、 $\mathbf{b} \in \mathbb{R}^d$  はバイアスである。また  $\mathbf{w}_i^{\text{dec}}$  は  $\mathbf{W}^{\text{dec}}$  の  $i$  列目のベクトルである。

本研究では JumpReLU SAE を研究の対象とし、JumpReLU SAE を SAE と略記する。

### 3 提案手法: SAE-MD

本研究の目的は SAE の“特徴”の学習を安定化させるために、SAE の学習の良い初期値を設定することである。通常の SAE の学習では  $\mathbf{W}^{\text{enc}}, \mathbf{W}^{\text{dec}}$  は一様分布で初期化される [8]。この初期化方法は勾配の消失・爆発への対策として有効であるが、SAE が獲得する“特徴”の安定化には寄与せず、異なる乱数シードで異なる“特徴”が学習される [9]。

本研究では SAE の学習に有効な初期値として Transformer 層の MLP のデコーダを用いることを提案する。以下では Transformer の内部構造の再解釈を通じてこのアイデアについて説明する。

**Transformer 層の再解釈** LLM の第  $l$  層への入力を  $\mathbf{x}^{(l)}$  とする。このとき、一般的な Pre-LayerNorm 型の Transformer 層の処理は

$$\mathbf{z}^{(l)} = \mathbf{x}^{(l)} + \text{MHA}(\text{LN}(\mathbf{x}^{(l)})), \quad (4)$$

$$\mathbf{x}^{(l+1)} = \mathbf{z}^{(l)} + \text{FFN}(\text{LN}(\mathbf{z}^{(l)})), \quad (5)$$

と記述される。ここで MHA, LN, FFN はそれぞれ Multi-Head の自己注意機構, レイヤー正規化, そして MLP を用いた非線形変換を表している。さらに式 (5) の第二項を書き下すと

$$\mathbf{U}^{\text{dec}} \sigma(\mathbf{U}^{\text{enc}} \text{LN}(\mathbf{z}^{(l)}) + \mathbf{c}^{\text{enc}}) + \mathbf{c}^{\text{dec}} = \sum_i \sigma_i \mathbf{u}_i^{\text{dec}} + \mathbf{c}^{\text{dec}},$$

となる。ただし  $\sigma$  は MLP の活性化関数で、 $\sigma_i$  は  $\sigma(\mathbf{U}^{\text{enc}} \text{LN}(\mathbf{z}^{(l)}) + \mathbf{c}^{\text{enc}})$  の第  $i$  成分である。また、 $\mathbf{U}^*$ ,  $\mathbf{c}^*$  は MLP のエンコーダおよびデコーダの重み行列とバイアスで、 $\mathbf{u}_i^{\text{dec}}$  はデコーダの重み行列  $\mathbf{U}^{\text{dec}}$  の  $i$  列目のベクトルである。式 (5) において第二項が支配的だと仮定できるとする<sup>1)</sup>と

$$\mathbf{x}^{(l+1)} \approx \sum_i \sigma_i \mathbf{u}_i^{\text{dec}} + \mathbf{c}^{\text{dec}}, \quad (6)$$

つまり LLM 内部の埋め込み表現  $\mathbf{x}^{(l+1)}$  は MLP のデコーダの重みのベクトル  $\mathbf{u}_i^{\text{dec}}$  の和で記述できる。

**本研究のアイデア**  $\mathbf{x}^{(l+1)}$  を SAE により分解・再構成すること、つまり式 (1), (2) において  $\mathbf{x} = \mathbf{x}^{(l+1)}$  とすることを考える。式 (6) と式 (2) を見比べると、下記の通り両者ともに  $\mathbf{x}^{(l+1)}$  を異なる方法でベクトルの和として表現していると解釈できる。

$$\mathbf{x}^{(l+1)} \approx \underbrace{\sum_i \sigma_i \mathbf{u}_i^{\text{dec}} + \mathbf{c}^{\text{dec}}}_{\text{MLP 近似}} \approx \underbrace{\sum_i z_i \mathbf{w}_i^{\text{dec}} + \mathbf{b}}_{\text{SAE 近似}}. \quad (7)$$

1) 式 (6) および以降では式 (5) の第一項  $\mathbf{z}^{(l)}$  を無視している。この項の適切な取り扱いは今後の課題である。

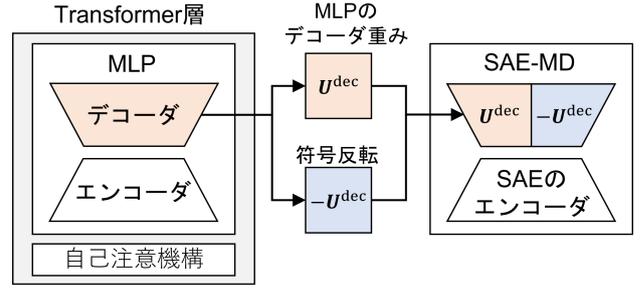


図 1: 提案手法 SAE-MD の概略図。LLM 内部の対応する Transformer 層の MLP のデコーダを SAE のデコーダの初期値として採用する。

この事実は SAE のデコーダのベクトル  $\mathbf{w}_i^{\text{dec}}$  として、MLP のデコーダのベクトル  $\mathbf{u}_i^{\text{dec}}$  が有効である可能性を示唆している。もちろん SAE は  $\mathbf{z}$  にスパース性を要求するため、 $\mathbf{u}_i^{\text{dec}}$  をそのまま SAE のデコーダに用いることが有効であるとは言い難い。しかし、 $\mathbf{u}_i^{\text{dec}}$  を  $\mathbf{w}_i^{\text{dec}}$  の初期値として用いることは有効であると考えられる。

**提案手法: SAE-MD** 上記のアイデアに基づいて、本研究では MLP のデコーダを SAE のデコーダの初期値に用いることを提案する。提案手法の概要を図 1 に示す。SAE では JumpReLU によって各  $z_i$  が非負であるのに対して、一般の MLP の活性化の値  $\sigma_i$  は負の値も取りえる。この場合、式 (7) の近似を成り立たせるには  $\sigma_i \mathbf{u}_i^{\text{dec}} = |\sigma_i| \times \text{sign}(\sigma_i) \mathbf{u}_i^{\text{dec}} \approx z_i \mathbf{w}_i^{\text{dec}}$  が必要であり、つまり  $\mathbf{u}_i^{\text{dec}}$  だけでなく  $-\mathbf{u}_i^{\text{dec}}$  も  $\mathbf{w}_i^{\text{dec}}$  の初期値の候補とする必要がある。そこで提案手法では  $\mathbf{U}^{\text{dec}}$  とその符号を反転させた  $-\mathbf{U}^{\text{dec}}$  を結合し、 $[\mathbf{U}^{\text{dec}}, -\mathbf{U}^{\text{dec}}] \in \mathbb{R}^{d \times 2m}$  を SAE のデコーダ  $\mathbf{W}^{\text{dec}}$  の初期値として用いる。ただし  $m$  は MLP の中間層の次元である。このようにデコーダを鏡状に構成することから、提案手法を Mirror Decoder (MD) を用いた SAE として SAE-MD と呼称する。

## 4 実験

提案手法 SAE-MD が (i) SAE の学習の安定化を実現すること、(ii) SAE としての品質が通常の SAE と同等であること、についての実験結果を報告する。

### 4.1 SAE の学習手順

LLM として Gemma-2-2B<sup>2)</sup> を採用し、第 12 層の埋め込み表現を使い通常の SAE および提案手法 SAE-MD による学習を行った。Gemma-2-2B の MLP

2) [hf.co/google/gemma-2-2b](https://hf.co/google/gemma-2-2b)

の中間層の次元は 9,216 のため、本実験の SAE および SAE-MD の中間層の次元は 18,432 とした。SAE の学習には `dictionary_learning` を用いた<sup>3)</sup>。SAE の学習データセットには `Pile-uncopyrighted`<sup>4)</sup> を使用した。ハイパーパラメータは付録 A に示す。

実験では SAE, SAE-MD の両者について中間層のスパース性を調整する損失関数を用いて異なるスパース性を有する複数の SAE (非ゼロ要素数 20, 40, 80, 160, 320) を学習した。また、それぞれの学習時の乱数シードは 24 および 42 に固定した二通りを実施した。これにより SAE の初期化および学習データのシャッフルが、SAE および SAE-MD の学習に与える影響を調査した。

## 4.2 学習の安定性の評価

異なる乱数シード 24 および 42 で学習された SAE (SAE<sub>24</sub>, SAE<sub>42</sub>) および SAE-MD (SAE-MD<sub>24</sub>, SAE-MD<sub>42</sub>) の間で学習された“特徴”がどの程度異なるのかを評価した。

安定性の評価は Paulo らの手法 [6] に従って、ハンガリー法を用いた“特徴”のマッチングを行った。以下では SAE<sub>24</sub>, SAE<sub>42</sub> を例に説明する。SAE-MD<sub>24</sub>, SAE-MD<sub>42</sub> についても同様の処理を適用する。

**1. 類似度行列の作成:** SAE<sub>24</sub>, SAE<sub>42</sub> についてコサイン類似度  $C_{ij}^{\text{enc}} = \cos(\mathbf{w}_{24,i}^{\text{enc}}, \mathbf{w}_{42,j}^{\text{enc}})$  および  $C_{ij}^{\text{dec}} = \cos(\mathbf{w}_{24,i}^{\text{dec}}, \mathbf{w}_{42,j}^{\text{dec}})$  を計算する。 $\mathbf{w}_{a,i}^{\text{enc}}, \mathbf{w}_{a,i}^{\text{dec}}$  は SAE<sub>a</sub> の  $\mathbf{W}^{\text{enc}}$  の  $i$  行目のベクトルおよび  $\mathbf{W}^{\text{dec}}$  の  $i$  列目のベクトルである。

**2. ハンガリー法の適用:** 類似度行列  $C^{\text{enc}}, C^{\text{dec}}$  のそれぞれに対して類似度を最大化するようにハンガリー法を使ってマッチングを行う。第  $i$  行にマッチングした列のインデックスをそれぞれ  $\pi(i), \sigma(i)$  とする。

上記の手順で“特徴”の比較を行った結果を図 2 に示す。図では SAE および SAE-MD のそれぞれで、各  $i$  について  $C_{i\pi(i)}^{\text{enc}}$  を縦軸、 $C_{i\sigma(i)}^{\text{dec}}$  を横軸にプロットし、ハンガリー法によりエンコーダとデコーダの両者で同じマッチングがされた場合 ( $\pi(i) = \sigma(i)$ ) に青、そうでない場合 ( $\pi(i) \neq \sigma(i)$ ) にオレンジで表記している。ハンガリー法で同じマッチングがされ、さらにコサイン類似度  $C_{i\pi(i)}^{\text{enc}}, C_{i\sigma(i)}^{\text{dec}}$  が 1.0 に近いときに、二つの異なる SAE で“同じ特徴”が学習されたと判定できる。図では、通常の SAE ではコサイン類似度が 0.6 以上のある程度類似した“特徴”が

3) [github.com/sapmarks/dictionary\\_learning](https://github.com/sapmarks/dictionary_learning)  
4) [hf.co/datasets/monology/pile-uncopyrighted](https://hf.co/datasets/monology/pile-uncopyrighted)

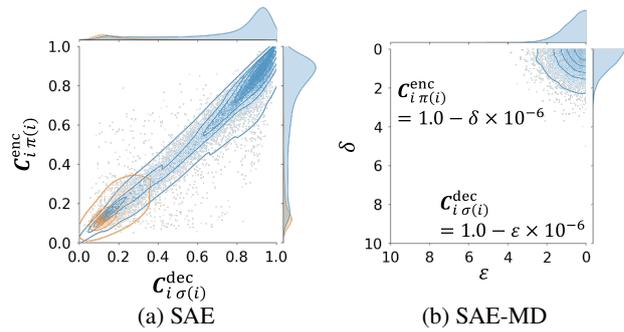


図 2: 異なる乱数シードで学習された SAE の特徴の一致度合い (中間層の非ゼロ要素数 320 の場合)。縦軸が  $C_{i\pi(i)}^{\text{enc}}$ 、横軸が  $C_{i\sigma(i)}^{\text{dec}}$  で、ハンガリー法によりエンコーダとデコーダの両者で同じマッチングがされた場合 ( $\pi(i) = \sigma(i)$ ) に青、そうでない場合 ( $\pi(i) \neq \sigma(i)$ ) にオレンジで表記している。

多く見られつつも、0.2 前後の小さな類似度が見られる“特徴”も一定数存在する。この結果は先行研究 [6] の報告とも整合的であり、SAE が乱数シードによって異なる“特徴”を無視できない割合学習していることを示している。他方、提案手法 SAE-MD ではコサイン類似度はほぼ 1.0 であり乱数シードが学習される“特徴”に影響を及ぼしていないこと、つまり SAE の学習の安定化に強く寄与していることが確認できる。

## 4.3 SAE としての品質の評価

前節では SAE-MD が学習の安定化に強く寄与することを確認した。本節では SAE-MD が従来の SAE と比較してその品質も同等であることを確認する。SAE の品質の包括的なベンチマークとして SAEBench を使用した [10]。SAEBench は、SAE の主要な能力である**再構成**、**解釈性**、**概念検出**、**因果的介入・分離**について 8 つの指標で評価を行う。下記に各指標の概要を示す。なお、どの指標ともに「大きい方が良い結果を示す」ように設計されている。

### 再構成

**Loss Recovered:** LLM の損失が、SAE による再構成後にどの程度復元されるかの割合。

### 解釈性

**Automated Interpretability:** SAE の各“特徴” $z_i$  の説明に基づき、未知のデータでの  $z_i$  の発火をどの程度正確に予測できるか。

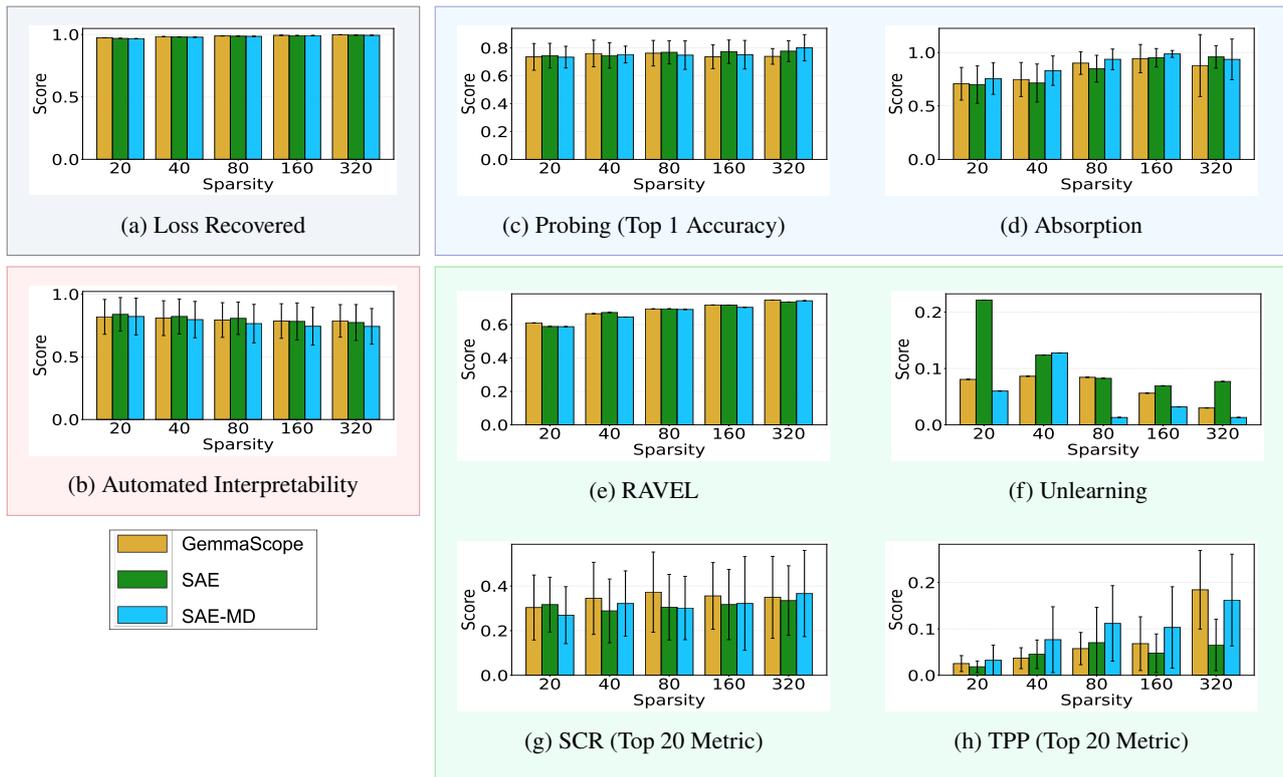


図 3: SAE Bench の結果 (平均と標準偏差)

### 概念検出

**$k$ -Sparse Probing:** 上位  $k$  個の  $z_i$  のみを用いた線形分類器による、特定概念の識別精度。

**Absorption:** 下位概念が上位概念の特徴量に「吸収」され、概念の識別が困難になっている度合い。

### 因果的介入・分離

**RAVEL:** 属性間の相関を排除し、特定の属性のみを独立して変更できる操作の正確性。

**Unlearning:** 標的とする知識の消去能力と、無関係なタスクにおける性能維持のバランス。

**Spurious Correlation Removal (SCR):** 介入によって「偽の相関」を除去し因果関係を抽出する能力。

**Targeted Probe Perturbation (TPP):** 多クラス設定において、特定クラスの特徴量への介入が他クラスの予測に与える影響 (因果的寄与度)。

図 3 に SAE Bench による評価の結果を示す。本実験で用いた SAE および SAE-MD の性能を客観的に評価するため、Gemma Scope<sup>5)</sup> [11] を比較対象として含めた。なお、乱数シードによる性能差はほぼ見られなかったため、ここでは乱数シードが 42 の結果を報告する。乱数シードの違いによ

る結果の差異は付録 B に示す。Loss Recovered と Automated Interpretability では SAE と SAE-MD がほぼ同等の値となっており、再構成と解釈性の面では両者は同等と言える。同様に、概念検出を測る指標 ( $k$ -Sparse Probing、Absorption) や、特徴の解離を測る指標 (RAVEL、SCR) についても、平均に小さな差はあれど標準偏差を考慮すると SAE と SAE-MD は同等の性能と言える。なお Unlearning と TPP は例外的にどの手法とも 0.1 前後とスコアが低く、SAE Bench [10] でのオリジナルの評価と整合的である。これらの結果は、提案手法 SAE-MD が安定性を大きく向上した上で、SAE の品質としては従来のものと比較して遜色ないことを示している。

## 5 おわりに

本研究では SAE の学習を安定化させる方法として、SAE の初期値に Transformer 層の MLP のデコーダを用いる手法 SAE-MD を提案した。実験の結果、SAE-MD は既存の SAE と同等の性能を有し、学習の安定性を向上できることを確認した。より詳細な実験評価および式 (6) で無視した式 (5) の第一項  $z^{(l)}$  の適切な取り扱いが今後の課題である。

5) [hf.co/google/gemma-scope-2b-pt-res](https://hf.co/google/gemma-scope-2b-pt-res)

## 謝辞

本研究は、JST 経済安全保障重要技術育成プログラム JPMJKP24C3 の支援を受けたものです。

## 参考文献

- [1] Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety – a review. **arXiv preprint arXiv:2404.14082**, 2024.
- [2] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. **ACM Transactions on Intelligent Systems and Technology**, Vol. 15, No. 2, 2024.
- [3] Andrew Y. Ng. Sparse autoencoder. **CS294A Lecture Notes, Stanford University**, pp. 251–258, 2011.
- [4] Robert Huben, Hoagy Cunningham, Logan Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In **International Conference on Learning Representations**, 2024.
- [5] Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In **International Conference on Representation Learning**, 2025.
- [6] Gonçalo Paulo and Nora Belrose. Sparse autoencoders trained on the same data learn different features. **arXiv preprint arXiv:2501.16615**, 2025.
- [7] Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. **Google DeepMind**, 2024.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, 2015.
- [9] Gonçalo Paulo, Alex Mallen, Caden Juang, and Nora Belrose. Automatically interpreting millions of features in large language models. In **International Conference on Representation Learning**, 2025.
- [10] Adam Karvonen, Can Rager, Johnny Lin, Curt Tigges, Joseph Bloom, David Chanin, Yeu-Tong Lau, Eoin Farrell, Callum McDougall, Kola Ayonrinde, et al. SAEBench: A comprehensive benchmark for sparse autoencoders in language model interpretability. In **International Conference on Machine Learning**, 2025.
- [11] Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. **arXiv preprint arXiv:2408.05147**, 2024.

## A SAE 学習のハイパーパラメータ

表 1: これらのハイパーパラメータは Karvonen ら [10] の設定を参考としている。

Hyperparameter	Value
Tokens processed	500M
Learning rate	$3 \times 10^{-4}$
Learning rate warmup (from 0)	1,000 steps
Sparsity penalty warmup (from 0)	5,000 steps
Learning rate decay (to 0)	Last 20% of training
Dataset	Pile-uncopyrighted
Batch size	2,048
LLM context length	1,024

## B 乱数シードによる SAEBench の結果の差異

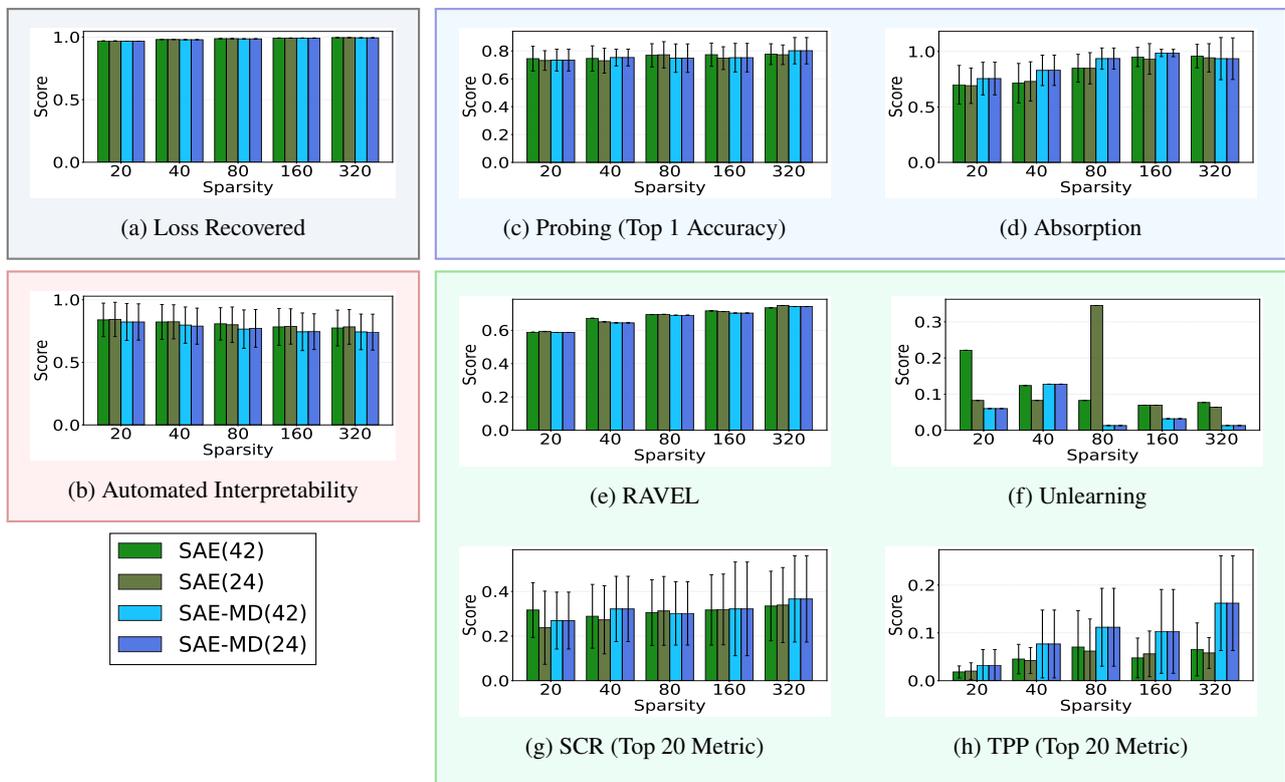


図 4: SAEBench の結果 (平均と標準偏差)。括弧内の数字は乱数シードの値を示す。