

“Task-relevant autoencoding” enhances machine learning for human neuroscience

Authors: Seyedmehdi Orouji¹, Vincent Taschereau-Dumouche²⁻³, Aurelio Cortese⁴, Brian Odegaard⁵, Cody Cushing⁶, Mouslim Cherkaoui⁶, Mitsuo Kawato⁴, Hakwan Lau⁷, & Megan A. K. Peters^{1,8}

1 Department of Cognitive Sciences, University of California, Irvine, Irvine, California, USA 92697

2 Department of Psychiatry and Addictology, Université de Montréal, Montreal, Canada, H3C 3J7.

3 Centre de recherche de l'institut universitaire en santé mentale de Montréal, Montréal, Canada.

4 ATR Computational Neuroscience Laboratories, Kyoto, Japan 619-0288

5 Department of Psychology, University of Florida, Gainesville, FL USA 32603

6 Department of Psychology, University of California Los Angeles, Los Angeles, 90095, USA

7 RIKEN Center for Brain Science, Tokyo, Japan

8 Center for the Neurobiology of Learning and Memory, University of California, Irvine, Irvine, California, USA 92697

Correspondence should be directed to:

Seyedmehdi Orouji
Department of Cognitive Sciences
2201 Social & Behavioral Sciences Gateway
University of California, Irvine
Irvine, CA 92697
sorouji@uci.edu

Megan A. K. Peters
Department of Cognitive Sciences
2201 Social & Behavioral Sciences Gateway
University of California, Irvine
Irvine, CA 92697
megan.peters@uci.edu

Abstract

In human neuroscience, machine learning can help reveal lower-dimensional neural representations relevant to subjects' behavior. However, state-of-the-art models typically require large datasets to train, so are prone to overfitting on human neuroimaging data that often possess few samples but many input dimensions. Here, we capitalized on the fact that the features we seek in human neuroscience are precisely those relevant to subjects' behavior. We thus developed a Task-Relevant Autoencoder via Classifier Enhancement (TRACE), and tested its ability to extract behaviorally-relevant, separable representations compared to a standard autoencoder for two severely truncated machine learning datasets. We then evaluated both models on fMRI data where subjects observed animals and objects. TRACE outperformed both the autoencoder and raw inputs nearly unilaterally, showing up to 30% increased classification accuracy and up to threefold improvement in discovering “cleaner”, task-relevant representations. These results showcase TRACE's potential for a wide variety of data related to human behavior.

Keywords: human neuroscience, machine learning, dimensionality reduction, task-relevant representation, fMRI, MVPA, autoencoder

1. Introduction

In studying the human brain and human behavior, we often use machine learning methods to home in on the (ideally lower-dimensional¹⁻⁴) representations contained in multivariate, feature-rich datasets. These data typically contain noisy, irrelevant signals¹⁹⁻²¹ that we would like to filter out using methods such as multivariate decoders⁵⁻⁸, various types of autoencoders, generative adversarial networks like InfoGAN⁹, or even principal components analysis (PCA)¹⁰⁻¹². However, state-of-the-art machine learning methods typically require very large datasets to train while data for individual human subjects collected with methods such as functional magnetic resonance imaging (fMRI)¹³⁻¹⁵ are often severely limited in sample size^{16,17} (i.e. have very few training exemplars compared to the dimension of data). Consequently, these methods are susceptible to overfitting on such neuroimaging data, reducing their predictive power and utility¹⁸⁻²⁰. What's more, parametric methods (such as PCA), which may better avoid the need for large training sets, by definition require rigid assumptions regarding the nature of the dimensionality reduction process and thus are limited *a priori* to insights consistent with these parametric assumptions. Thus, we are in need of a nonparametric method that can reveal the *low-dimensional, task-relevant* representations in a given brain region using *exemplar-poor but input-dimension-rich* datasets.

Here, we sought to capitalize on a unique property of many human neuroimaging datasets, which is that the features we wish to identify can be conceptualized based on whether they are relevant for the subject's behavior.

We drew inspiration from previous successes with classifier-enhanced autoencoders²¹⁻²⁴ to develop the Task-Relevant Autoencoder via Classifier Enhancement (TRACE) model. TRACE's architecture is purposely simple to limit overfitting to small datasets, consisting of a fully-connected autoencoder with only one hidden layer on each of the encoding and decoding arms and a logistic regression classifier attached to the bottleneck layer (**Figure 1**).

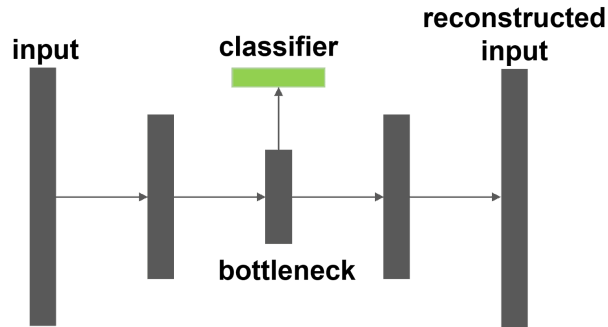


Figure 1. A cartoon representation of the TRACE network architecture. Input is connected to the bottleneck via one hidden encoding layer, and again to the reconstruction via one hidden decoding layer. A classifier is attached to the bottleneck and contributes to the objective optimization function.

We developed four quantitative metrics to assess TRACE's performance at different bottleneck dimensionalities (compression levels), and then comprehensively benchmarked TRACE under conditions of severe data sparsity using the MNIST²⁵ and Fashion MNIST²⁶ datasets, two of the most popular machine learning datasets. We then applied TRACE to a neuroimaging (fMRI) dataset of subjects who viewed and categorized animals and objects while blood oxygen level dependent (BOLD) signal was collected from ventral temporal cortex (VTC) in a single, 1-hour

session. By constraining the dimensionality reduction process to specifically prioritize features that were relevant to the participants’ behavioral task, we show that TRACE can extract both quantitatively and qualitatively ‘cleaner’ representations at both reduced dimensions and in the original input dimensionality, showing up to threefold improvement in decoding accuracy and separation of class-specific patterns. These results demonstrate our method can distill highly separable, low dimensional neural representations even with sparse and noisy data. TRACE may thus show promise on a broad variety of behaviorally-relevant neuroimaging datasets.

2. Results

We quantified the performance of the Task-Relevant Autoencoder via Classifier Enhancement (TRACE) model against that of a standard AE with otherwise identical architecture via (1) *reconstruction fidelity*, (2) *reconstruction classifier accuracy*, (3) *bottleneck classifier accuracy*, and (4) *reconstruction class specificity* (see **Methods Section 4.4**) (“class” here refers to the class of the input image, e.g. “9” or “shoe” or “cat”). We assessed these metrics as a function of different bottleneck dimensionalities (i.e., compression levels), first on the MNIST and Fashion MNIST datasets under increasing data sparsity and then on a previously-collected fMRI dataset of ventral temporal cortex (VTC) (i.e., voxel activations while 59 human subjects viewed 40 classes of animals and objects). We also performed additional investigation at each dataset’s ‘optimal’ bottleneck dimensionality (where reconstruction class specificity is maximized) to characterize each model’s behavior.

2.1 Benchmarking TRACE's advantages, including under increasing data sparsity

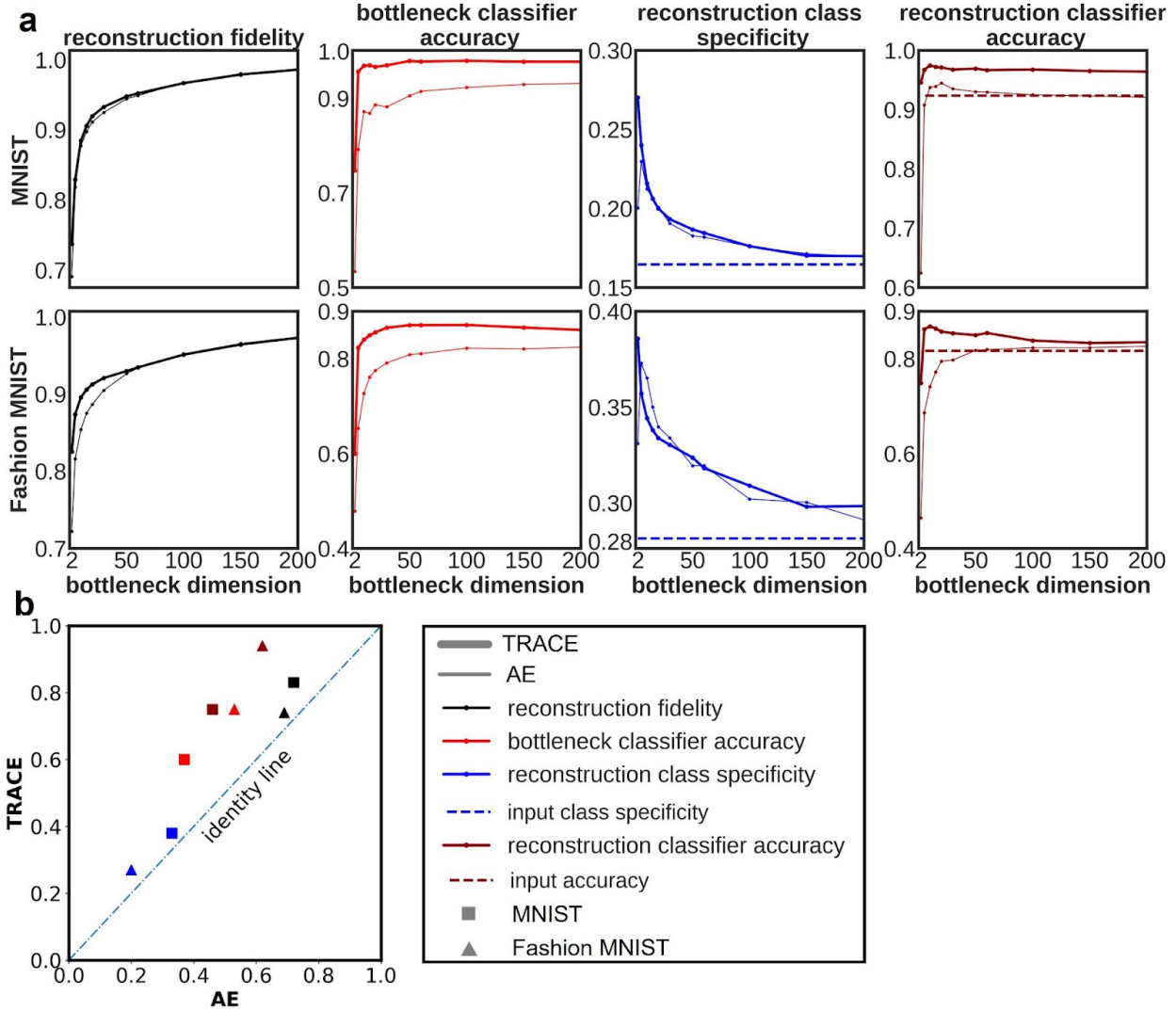


Figure 2. Quantitative comparison between TRACE and AE on the four outcome metrics, for the two benchmark datasets (MNIST & Fashion MNIST) for bottleneck dimensionalities between 2 and 200. All metrics show superiority of TRACE over AE. **(a)** TRACE is shown by the thicker line while AE is shown by thinner line. The black, dark red, red, and blue lines show the reconstruction fidelity ($fidelity_R$), reconstruction classifier accuracy (RCA), bottleneck classifier accuracy (BCA), and reconstruction class specificity (RCS), respectively (see **Methods**). The dashed dark red and blue lines show the input classifier accuracy and input class specificity, respectively. AE outcome metrics for all bottleneck dimensionalities tested (2-1500) are shown in **Figure S3**; locations of peaks for all four metrics are shown in **Table S1**. **(b)** Quantitative outcome metrics at the optimal bottleneck dimensionality value ($d=2$) for reconstruction class specificity, for the MNIST and Fashion MNIST datasets. See main text for details.

We first examined *reconstruction fidelity* (black, **Figure 2a**), i.e. the mean Pearson correlation of the inputs and corresponding reconstructions. High reconstruction fidelity assures us that the discovered features in the bottleneck provide a reasonable representation of this high-dimensional information – i.e., that the autoencoder portion of the models can be successfully trained. At lower bottleneck dimensionalities ($d \leq 30$), reconstruction fidelity was slightly higher for TRACE, becoming asymptotically equivalent between the models as bottleneck dimensionality increases ($d > 100$). Notably, TRACE’s superior reconstruction fidelity occurred despite the fact that the contribution of the reconstruction part of the loss function (mean square error; MSE) for TRACE was smaller than for AE (i.e., the objective function in TRACE is the sum of reconstruction loss (L_R) and classification loss functions (L_{CE}), but in AE reconstruction is the only objective; see **Methods Section 4.4.1**).

We next examined *reconstruction classifier accuracy* (dark red, **Figure 2a**), i.e. the accuracy of a separate logistic classifier trained to discriminate classes using reconstructed data. Reconstruction classifier accuracy quantifies the task-relevance of the information extracted through the compression process, and also provides a benchmark against which to compare directly to the original input (see below). Reconstruction classifier accuracy for MNIST peaked at bottleneck dimensionality $d=20$ for AE and $d=10$ for TRACE, and was consistently higher for TRACE over AE. A similar pattern was observed for Fashion MNIST, with a peak at $d=10$ for TRACE but an ever-increasing monotonic function with no peak for AE, with TRACE always higher than AE. Interestingly, that this metric peaks at higher bottleneck dimensionalities than other metrics suggests that the performance of a classifier trained on these high-dimensional reconstructions may not meaningfully reflect the maximum compression that TRACE can achieve without loss of overall performance. Below, we also compare reconstruction classifier accuracy to input classifier accuracy as a measure of TRACE vs AE’s capacity to extract task-relevant information.

Next, we examined *bottleneck classifier accuracy* (bright red, **Figure 2a**), i.e. the accuracy of a separate classifier trained with bottleneck features as input after the training of the main AE or TRACE model. Bottleneck classifier accuracy was much higher for TRACE than for AE even at very low bottleneck dimensionalities. As bottleneck dimensionality grew, this metric asymptotically equalizes to around 5% and 4% better in the MNIST and Fashion MNIST datasets, respectively. Notably, though, in both datasets, at all bottleneck dimensionalities tested, TRACE bottleneck classifier accuracy is *always* higher than that for AE.

The fourth metric we examined was *reconstruction class specificity* (blue, **Figure 2a**), i.e. the average within-class correlation of the reconstructed inputs minus the average between-class correlation. This metric quantifies the degree of separation between class clusters in reconstruction feature space as a measure of reconstruction representations’ categorical ‘purity’. Reconstruction class specificity peaks at bottleneck dimensionality $d=2$ for TRACE for both MNIST and Fashion MNIST. As with the other metrics, TRACE outperformed AE in almost all cases, with brief exceptions between $d=5$ and $d=20$ (and $d > 1000$; see **Supplementary Material Figure S2**) for Fashion MNIST.

A final – and critical – test of TRACE would examine its ability to not only distill task-relevant information into low-dimensional representations but also ‘push’ such distilled insights back into the native space of the input. This would be especially important if one wished to use TRACE to

de-noise fMRI data to discover multivoxel patterns representing a target concept or category to be used with noninvasive intervention strategies such as DecNef^{27–30}. Although iterative sparse logistic regression and support vector machine classification have been demonstrated as successful at identifying such patterns when trained on the native input data^{27,31,32}, we wanted to see whether TRACE would be able to denoise the data such that an even cleaner target pattern would become discoverable. Specifically, if TRACE is successful at actively removing task-irrelevant noise rather than simply passively averaging across it (as is done with a standard category-based logistic regression) or removing it through iterative sparsity approaches (iterative sparse logistic regression), then we should observe two patterns. First, reconstruction classifier accuracy should approach or exceed classification accuracy of an identical logistic regression classifier trained on the native inputs. Second, reconstruction class specificity should behave similarly, approaching and then exceeding input class specificity. This behavior makes reconstruction class specificity an ideal metric for defining the ‘optimal bottleneck dimensionality’ if one’s goal is to optimally distill representations in native space.

To evaluate this behavior, we (a) trained an additional logistic regression classifier on each of the datasets to classify the native input, and (b) computed class specificity directly from the raw input data for all three datasets. We then compared the outcomes to the reconstruction classifier accuracy and reconstruction class specificity computed as a function of bottleneck dimensionality.

Results revealed that, for MNIST, reconstruction classifier accuracy (solid dark red, **Figure 2a**) exceeded input classifier accuracy (dashed dark red line) immediately (at $d=2$) for TRACE but not until $d=10$ for AE (and then only until $d=20$, at which point it falls again). For Fashion MNIST, this occurred at $d=5$ for TRACE and $d=100$ for AE. These results show that TRACE provides not only superior compression but also superior denoising even in comparison to the direct inputs, over the standard AE model. TRACE’s denoising capability can be particularly useful in DecNef^{27–29,33–39} studies as it can minimize the task-irrelevant information of exemplars even in the anatomical and functional brain space.

Results for reconstruction class specificity followed a different pattern, but still favored TRACE: reconstruction class specificity (solid blue line, **Figure 2a**) exceeded input class specificity (dashed blue line) at most bottleneck dimensionalities for both TRACE and AE, but was higher for TRACE than AE. These results show that TRACE can provide a powerful method for not only distilling low-dimensional representations, but also in pushing those cleaner representations back into the structure and dimensionality of the raw input space. That is, a structurally identical logistic classifier with the same number of parameters can exhibit better performance using the reconstructed inputs than using the original inputs themselves.

2.1.1 Comprehensive comparison across metrics as a function of increasing data sparsity

We next sought to select a single bottleneck dimensionality for TRACE to explore its benefits over AE under increasing data sparsity. For this purpose, we selected the maximal value of reconstruction class specificity because this metric provides the best balance between task-relevant information extraction and compression, both for analyzing low-dimensional representations and patterns in the original input dimensionality (e.g. for use with real-time DecNef^{27–29,33–39}).

Reconstruction class specificity peaked at $d=2$ for both MNIST and Fashion MNIST, so we can first examine TRACE's superiority at this dimensionality when maximal data is available ($n = 60,000$ training samples for both datasets). Here, we see that TRACE's superior extraction of task-relevant information comes at no loss in reconstruction fidelity over AE (**Figure 2b**; **Table S2**). Further explorations, described below, were therefore done at bottleneck dimensionality $d=2$.

To examine how TRACE versus AE fare under increasing data sparsity, we trained each model after removing 10, 30, 50, 70, 90, 95, and 98 percent of the training data. Training examples at each level of sparsity for both TRACE and AE remained the same. We then used the conventional 10,000 held-out test set on the trained models and calculated all four metrics for all levels of data sparsity.

TRACE was much more robust to increasing data sparsity than AE (**Figure 3**). Specifically, AE showed immediate drops in all four metrics as available training data decreased, whereas TRACE's performance was strikingly stable until only 2% of the data (1200 samples) remained available for training. We note that the fMRI dataset we use below has a similar samples-to-input-dimensions ratio as the 98% truncated MNIST and Fashion MNIST datasets (~ 1.6 for MNIST and Fashion MNIST, and ~ 1.5 for this fMRI dataset). At this level of data reduction (i.e. 98% truncation) and bottleneck dimensionality $d=2$, we performed 1000 bootstrap replications to randomly remove 98 percent of exemplars in MNIST and Fashion MNIST and reported the mean values of the 1000 replications for all metrics. As shown in **Figure 3**, TRACE continued to demonstrate superior performance even at the most extreme level of data truncation (i.e. 98% truncation). TRACE nearly uniformly swept AE across all performance metrics. The probability of it outperforming AE out of 1000 replications was as follows. 1. Reconstruction fidelity: $p = 0.991$, $p = 0.55$ in MNIST and Fashion MNIST respectively; 2. Bottleneck classifier accuracy: $p = 0.755$, $p = 0.594$ in MNIST and Fashion MNIST respectively; 3. Reconstruction specificity: $p = 1$, $p = 0.921$ in MNIST and Fashion MNIST respectively; and 4. Reconstruction classification accuracy: $p = 1$, $p = 0.958$ in MNIST and Fashion MNIST respectively.

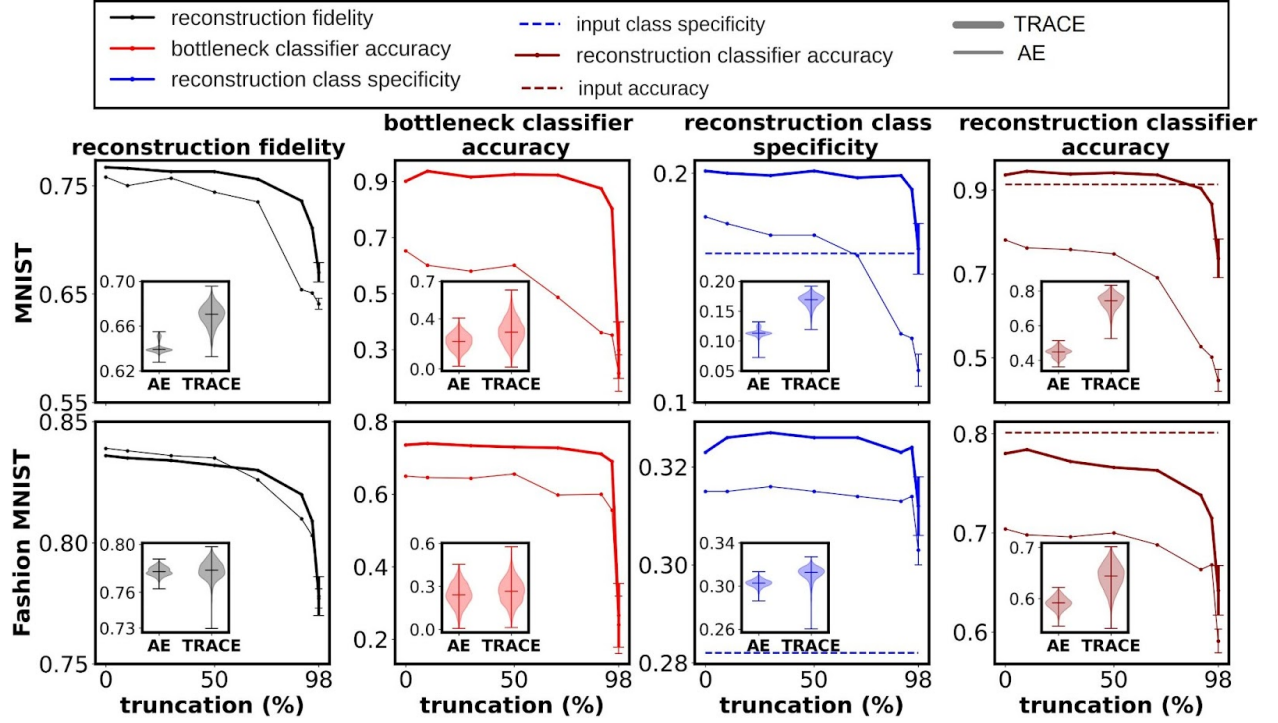


Figure 3. Performance of AE and TRACE model as a function of sample size for the optimal bottleneck dimension of $d=2$. At 98% truncation level, we used 1000 bootstrap replications to randomly truncate 98 percent of exemplars and reported the means and standard deviations of the metrics for MNIST and Fashion MNIST. Error bars show the standard deviation of results across 1000 bootstrap replications at 98% data truncation. The embedded plots show the distributions of 1000 bootstrap replications in AE and TRACE at 98% data truncation, with horizontal lines showing the medians of the distributions. The small variation in the metrics is likely due to random initialization of weights and use of GPUs in fitting the models.

At maximal data reduction (98% truncation) and bottleneck dimensionality $d=2$, we then performed additional explorations of both bottleneck representations and reconstructions. First, we visualized bottleneck representations by plotting the activities of the two bottleneck features against each other for each of the 10 classes in each dataset for TRACE versus AE (**Figure 4**). The results are striking: TRACE showed superior task-relevant representations especially for MNIST, i.e. a clear qualitative advantage in clustering performance showing distinct clusters for different classes (**Figure 4b**) in stark contrast to AE’s class clusters, which are heavily overlapping (**Figure 4a**). Although this difference in clustering ability was less apparent for Fashion MNIST (**Figure 4c,d**), TRACE’s clusters do appear visually more tightly bound.

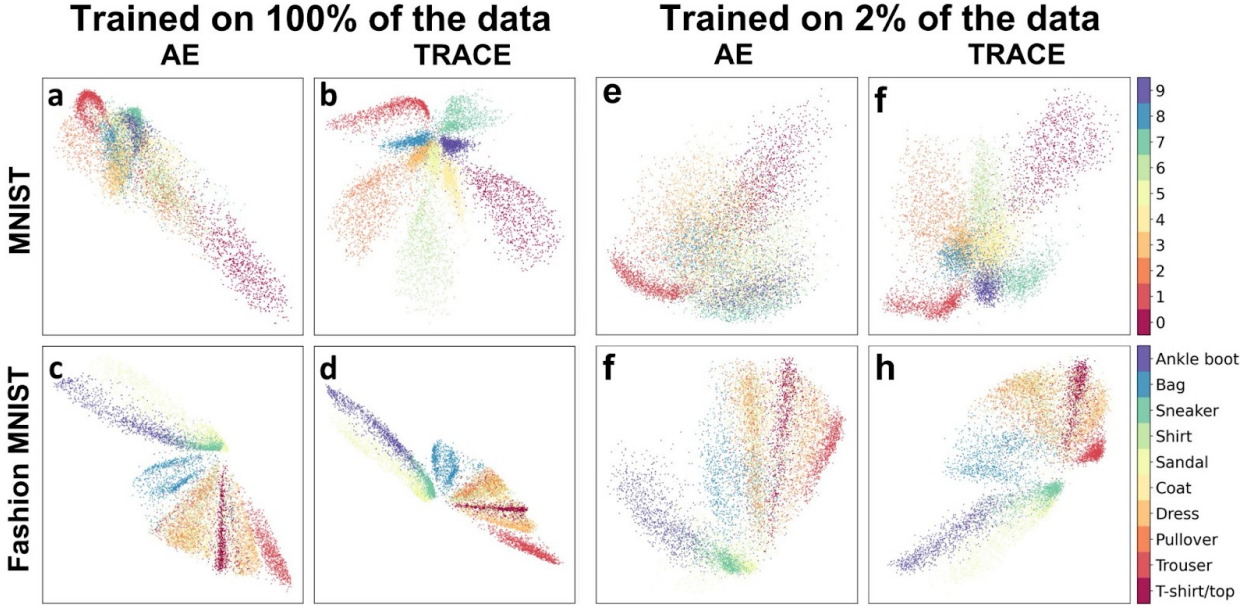


Figure 4. Visualization of bottleneck features for MNIST and Fashion MNIST datasets using AE and TRACE. **(a-d)** When trained on the full dataset, TRACE shows clear superiority in creating distinctive clusters in the bottleneck for different classes for MNIST dataset in comparison to AE. The distinction is less clear but still apparent in the Fashion MNIST dataset. **(e-h)** This pattern persists even when trained on only 2% of the dataset, again showing the robustness of TRACE.

We next turned to examining the reconstructions (still at bottleneck $d=2$). We first examined the MNIST reconstructions for several different exemplars of the same categories (e.g., several different “3” and “6” exemplars). TRACE’s superiority is clear to the naked eye: the reconstructions of particular “3” and “6” exemplars from TRACE are much more “three-like” and “six-like” than reconstructions from AE (**Figure 5a**). (Recall that this qualitative superiority does not come at any quantitative cost to the reconstructions; in fact, reconstruction fidelity was higher for TRACE than for AE at $d=2$, **Figure 2b**.) Similar findings held for Fashion MNIST (e.g. sandal and shirt, **Figure 5b**), although the visual result is less striking. These patterns held even when only 2% of the data was available for training (**Figure 5c & 5d**).

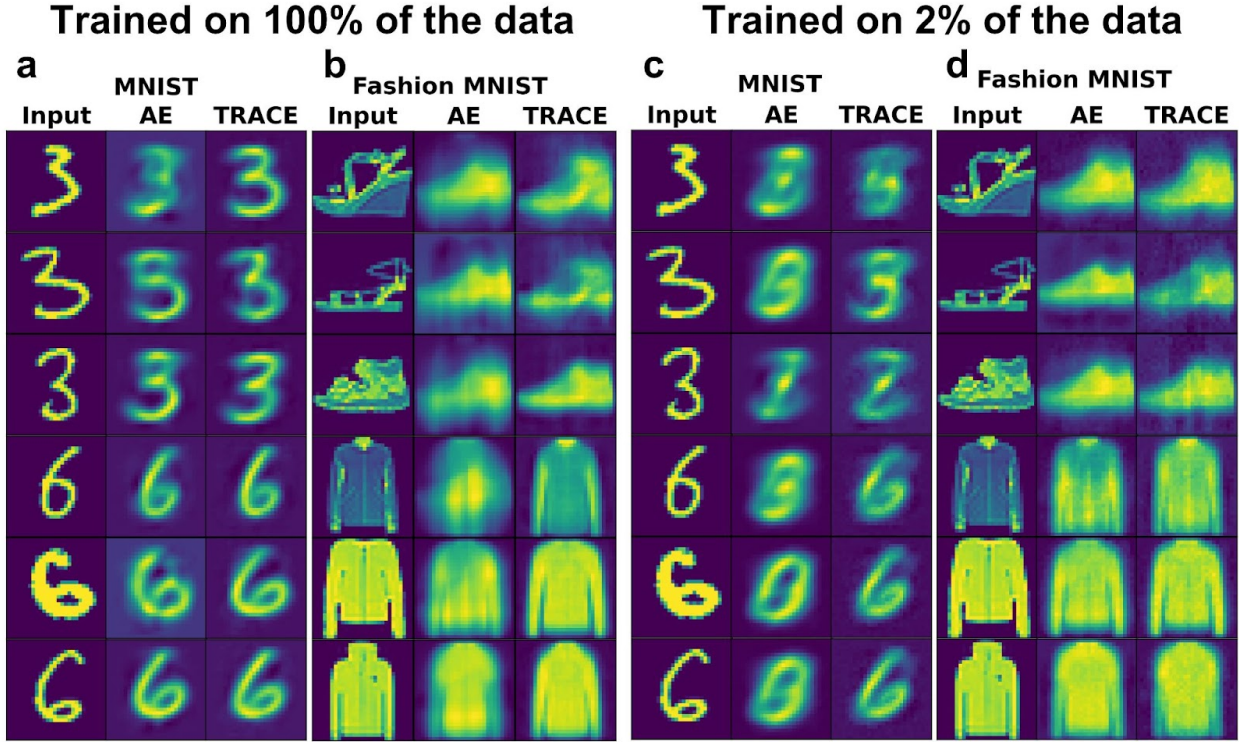


Figure 5. Visualization of reconstructions for MNIST and Fashion MNIST datasets using AE and TRACE. The reconstruction of three representative instances of numbers “three” and “six” in MNIST dataset and three instances of classes “sandal” and “shirt” in the fashion MNIST dataset, using AE, and TRACE networks when there are two features in the bottleneck shows the same pattern. TRACE shows a more clear and *canonical* reconstruction of the inputs across several exemplars from the same category.

We next wanted to quantitatively investigate the distributions of within-class versus between-class clusters, both in the bottleneck and the reconstructions. This approach will facilitate evaluation of the fMRI dataset since visual inspection in fMRI data is not possible in the same sense as for MNIST and Fashion MNIST given that optimal bottleneck dimensionality is larger than 2 (see **Results Section 2.2**). We computed the effect size (Cohen’s d) separating clusters in both the bottleneck and reconstructions using pairwise within- versus between-class Euclidean distances. Whether trained on all of the data or 98% truncated, Cohen’s d was always larger for TRACE than for AE (**Table 1**).

		MNIST	Fashion MNIST	MNIST (98% truncated)	Fashion MNIST (98% truncated)
Bottleneck	AE	1.07 ± 0.34	1.41 ± 0.46	1.01 ± 0.44	1.41 ± 0.41
	TRACE	1.7 ± 0.2	1.62 ± 0.52	1.39 ± 0.29	1.56 ± 0.39
Reconstruction	AE	1.31 ± 0.27	1.45 ± 0.61	1.02 ± 0.45	1.41 ± 0.65
	TRACE	1.6 ± 0.22	1.59 ± 0.69	1.41 ± 0.26	1.56 ± 0.76

Table 1. Cohen's d measures of effect size comparing within-class versus between-class Euclidean distances in the bottleneck and reconstructions for TRACE and AE.

2.2 TRACE's performance on a real fMRI dataset

Given TRACE's apparent superiority over AE even under extreme data sparsity, we next sought to evaluate TRACE using a real-world fMRI dataset, since ultimately our goal is to learn about neural representations. Thus, we used the same metrics as we used to evaluate TRACE on MNIST and Fashion MNIST on an fMRI dataset consisting of 59 individuals who each viewed 3600 exemplars of 40 classes of animals and objects (90 exemplars per class) while BOLD signal from ventral temporal cortex (VTC) was obtained.

Excitingly, the fMRI dataset showed the same patterns in our four quantitative metrics as the MNIST and Fashion MNIST datasets almost across the board. First, although reconstruction fidelity was slightly lower for TRACE than for AE for lower dimensions in the bottleneck ($d \leq 100$), it was slightly better for TRACE for higher dimensions ($d > 150$). This pattern is possibly due to the fact that the actual dimensionality of the low-dimensional representations in the fMRI dataset are likely to be higher than those of MNIST and Fashion MNIST, given that this dataset contains ventral temporal cortical representations – a high level ventral stream visual area encoding for complex and semantic features in objects.

Reconstruction classifier accuracy followed an ever-increasing value for both TRACE and AE, but again TRACE showed higher reconstruction classifier accuracy than AE at all bottleneck dimensionalities tested. TRACE also showed higher bottleneck classifier accuracy at all bottleneck dimensionalities in comparison to AE. Reconstruction classifier accuracy even surpassed the input classifier accuracy for bottleneck dimensionalities higher than $d=60$ (dashed blue line, **Figure 6a**) which again suggests that the reconstructed version in the original input space contains more task-relevant information.

TRACE outperformed AE in reconstruction class specificity as well, showing that even in the native space of the input – i.e., voxel patterns of activity in ventral temporal cortex – TRACE not only successfully distills lower-dimensional representations of task-relevant data, but also faithfully projects them back into original, high-dimensional voxel space. Reconstruction class specificity peaked at bottleneck dimensionality $d=30$, and then fell again. The same was not true for AE, for which reconstruction class specificity rose but then asymptoted. Crucially, though, reconstruction class specificity was also always higher for TRACE than for AE, much exceeding input class specificity (**Figure 6a**, solid and dashed blue lines, respectively). This capacity to distill a task-relevant, low-dimensional representation and put it back in brain space could potentially have great value for studies in which such multivoxel patterns are the target of DecNef^{27–30} or other investigations which require anatomically-related representations. We discuss this possibility in greater detail in the **Discussion**, below.

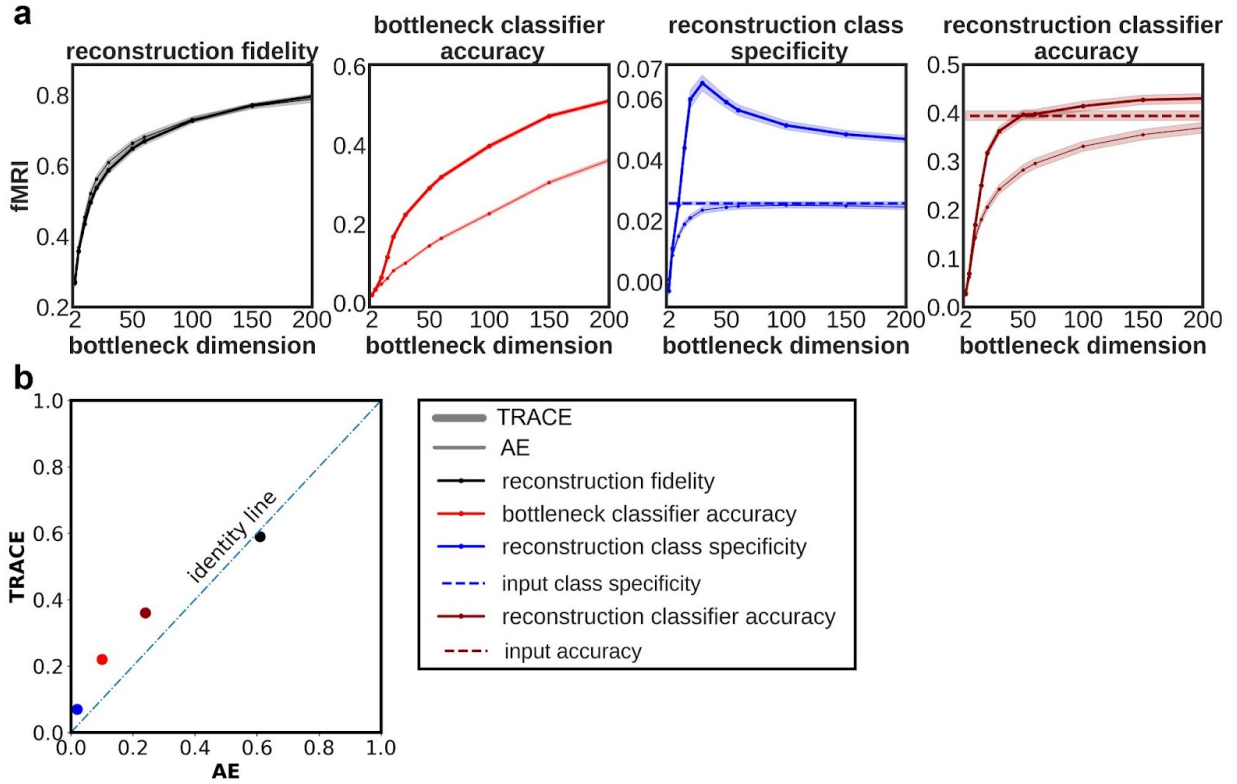


Figure 6. Comparison between quantitative metrics for AE and TRACE for the fMRI dataset ($n=59$). **(a)** TRACE shows superior performance in three out of four metrics (excluding reconstruction fidelity), especially reconstruction class specificity which peaks at $d=30$ features in the bottleneck. **(b)** Comparison of performance TRACE vs AE across the four at the optimal bottleneck dimensionality ($d=30$). Error bars (standard error) are so small as to be invisible. See main text for details.

2.2.1 Exploration at optimal bottleneck dimensionality for fMRI data

As mentioned above, the maximal value for reconstruction class specificity was found at $d=2$ for the MNIST and Fashion MNIST datasets. For the fMRI dataset, we found that reconstruction class specificity peaked at $d=30$, so we proceeded with a parallel analysis to that done above at this dimensionality.

Crucially, at $d=30$, TRACE's performance on the fMRI dataset mimicked its exemplary performance on the MNIST and Fashion MNIST datasets with the exception of reconstruction fidelity, which was only slightly smaller for TRACE than for AE (**Figure 6b**). Although a paired t-test on reconstruction fidelity did demonstrate that reconstruction fidelity for TRACE was statistically lower than that of AE (paired $t(58) = 4.6$, $p = 2e-5$), the effect size (i.e. Cohen's d) was 0.27 which is considered small. Recall, however, that reconstruction fidelity is a measure which is actually at odds with the goal of discovering canonical representations (either in the original input space or the bottleneck), since it captures both how well the model can reconstruct signal *and* how well it can reconstruct exemplar-specific noise. Importantly, the other metrics, which quantify the efficiency of task-relevant information distillation, show that TRACE performed superiorly to AE even in this fMRI dataset. TRACE outperformed AE in bottleneck classifier accuracy (paired $t(58) = 23.76$, $p = 2.9e-31$, Cohen's $d = 3.32$), reconstruction classifier accuracy ($t(58) = 12.1$, $p = 2.5e-17$, Cohen's $d = 0.91$), and reconstruction class

specificity ($t(58) = 23.7$, $p = 3.3e-31$, Cohen's $d = 2.8$). The effect size for these three metrics is considered large.

Ultimately, as our goal is to learn about representations in human VTC, we also might want to visualize clusters for the 40 classes of the fMRI dataset. However, unlike for MNIST and Fashion MNIST where optimal bottleneck dimensionality was $d=2$, for the fMRI dataset we found the optimal bottleneck dimensionality at $d=30$. Therefore, we cannot easily visualize the class clusters in a scatterplot, and performing further dimensionality reduction for the sake of visualization would be inappropriate since assumptions of whichever dimensionality reduction technique we chose would impact the visualizations. Instead, we can use the same Cohen's d approach, described above, to characterize the tightness of the class clusters even in higher dimensionalities. The average effect size separating within- and between-class Euclidean distances across all 59 subjects was $0.35 (\pm 0.09)$ for TRACE but only $0.14 (\pm 0.04)$ for AE, again showing TRACE's superiority.

As a final evaluation of TRACE's ability to filter out task-irrelevant information, we calculated the within- versus between-class Euclidean distance Cohen's d in the reconstructions. Pushing the distilled representations back into input space is particularly exciting for the use of TRACE with fMRI data if one wishes to discover a particular target pattern for further anatomical analysis, or for use with real-time neuroimaging (e.g., DecNef). However, visually examining fMRI reconstructions would not provide particularly useful information about the 'cleanliness' of the reconstruction, as the patterns are not visually meaningful to begin with, so we must again rely on a quantitative comparison. The average Cohen's d here again showed TRACE's superiority, with mean Cohen's d of $0.16 (\pm 0.05)$ across subjects for TRACE but only $0.1 (\pm 0.03)$ for AE. In other words, TRACE was able to reduce task-irrelevant information and thus extract a 'cleaner' representation, even in the reconstructions.

3. Discussion

3.1 Summary of findings

Most dimensionality-reduction approaches do not have a specific mechanism to ensure that the lower dimensional representations they reveal are particularly relevant to the question an experimenter wishes to answer. Further, many state-of-the-art deep learning models are of limited utility for discovering and characterizing meaningful representations in input-dimension-rich but exemplar-poor datasets, as they tend to overfit^{18–20}. Together, these facts make discovering neural representations in within-subject fMRI datasets – which also often contain a high degree of noise and task-irrelevant information – extremely challenging^{13–15}. Further, to address these issues we proposed the Autoencoder with Classifier Enhancement (TRACE) model: a simple autoencoder with a classifier attached to the bottleneck. The classifier forces the model to learn not just lower dimensional representations of the data, but those that are also task-relevant. To quantify TRACE's superiority over a standard autoencoder (AE), we used four metrics (see **Methods Section 4.4**): 1. reconstruction fidelity; 2. bottleneck classifier accuracy; 3. reconstruction class specificity; and 4. reconstruction classifier accuracy.

TRACE outperformed AE in all metrics, with the exception of reconstruction fidelity (sometimes). Moreover, at the 'optimal' bottleneck dimensionality, TRACE's superior capacity for extracting task-relevant information is evident in both the bottleneck and reconstruction, and TRACE's reconstructions can even outperform the inputs on a measure of task-relevant behavior (reconstruction class specificity). TRACE's advantage over AE appears due to TRACE's

capacity to minimize task-irrelevant, idiosyncratic information unique to a particular sample. This is evident in the one occasional exception to TRACE’s sweeping superiority: reconstruction fidelity for the fMRI dataset. However, this seeming underperformance – especially in the fMRI dataset – is actually a strength: AE tried “too hard” to encode idiosyncratic details of a particular exemplar in the bottleneck, when some of these details are merely noise for the task that the observer is performing. Thus, precise reconstruction of noisy data may not be suitable.

Critically, all of these behaviors were maintained by TRACE even under extreme data truncation for the MNIST and Fashion MNIST datasets, and carried over into a real-world fMRI dataset. These results suggest that TRACE can extract lower-dimensional representations of data for both reconstruction and classification purposes and can do so even when there is a highly undesirable balance of input-dimensions versus samples, as in typical fMRI data – suggesting strong promise for TRACE and variants of this approach for both fMRI datasets and for other biological-scale data with many more input-dimensions than samples.

3.2 Relation to previous work

TRACE is not the only model which can accomplish dimensionality reduction, but one of many techniques. So is TRACE really necessary? Why would principal components analysis (PCA)^{10–12} not suffice? PCA focuses on creating new features that can best explain the variance in data – including the noise and task-irrelevant information, which we know to be problematic especially in fMRI data^{14,40–44} – and thus lacks explicit mechanisms to ensure the discovered lower representations contain task-relevant information. Additionally, we also note that PCA-based methods are not assumption-free (that is, they are parametric); these assumptions about the functional form of the dimensionality reduction limit the discovered features to adhering to those assumptions.

Other techniques have been developed including nonparametric techniques beyond the fully-connected AE used here^{45,46}: variational (stochastic) autoencoders (VAEs)^{47,48}, adversarial autoencoders⁴⁹, generative adversarial networks (GANs)⁵⁰, deep convolutional GANs (DCGANs)⁵¹, and so on. While comprehensive exploration of these is beyond the scope of this manuscript, we note that many of these models do still suffer from the fact that the discovered lower dimensional representations are not explicitly crafted to be task-relevant⁵². This consideration led to the development of InfoGAN⁹, an unsupervised learning technique which modifies a generative adversarial network (GAN) in order to learn interpretable, low-dimensional representations. InfoGAN accomplishes this task by maximizing mutual information between noise in the GAN network and observations. Yet despite the tremendous success of InfoGAN⁹, it is highly disadvantaged for the limited (sample-poor) data type targeted here. Specifically, InfoGAN’s success has been demonstrated only on large-scale training datasets consisting of tens of thousands of training images.

Attempts to mitigate the curse of dimensionality in fMRI datasets by pooling data across subjects to create larger training sets have of course been established to try to mitigate this significant challenge, including the shared response model⁴, hyperalignment^{53–55}, and more recently decoder + autoencoder approaches⁵⁶. However, while these can pool fMRI data to create more training exemplars, they do not explicitly seek subject-specific response patterns and instead presuppose that all subjects share a common response pattern.

In sum, although we do not benchmark TRACE against InfoGAN, hyperaligned data, or the expansive space of model variants, we argue that TRACE’s utility is not only in its ability to distill task-relevant, low-dimensional representations, but also to do so in exemplar-limited,

biological-scale datasets such as those collected in human neuroimaging experiments within a single subject.

3.3 Limitations

One limitation of the present approach is that we (deliberately) made TRACE and AE extremely simple, which could have limited their performance. We did not investigate whether TRACE-like architecture (addition of a classifier on the bottleneck layer) would similarly improve performance for a variational autoencoder (VAE) network, or whether multi-layer perceptrons or convolutional neural network (CNN) classifiers would surpass the simple logistic regression classifiers used here. We also could have opted to make the models deeper, with many hidden layers, which might have resulted in benefits in classification or reconstruction. However, we reiterate that we selected a simple architecture to be able to best evaluate TRACE's advantages over a "plain vanilla" fully-connected autoencoder, as more complex architectures could obscure TRACE's advantages. Future work may wish to explore other possible TRACE-like architectures.

It is also worth mentioning that for the sake of consistency we kept all hyperparameters for all networks and datasets the same. However, during training TRACE on a new dataset, it is always possible to tune the hyperparameters (learning rate, batchsize, regularization, etc) in order to achieve better performance (e.g. better bottleneck classification accuracy). Future studies may also more comprehensively explore the impact of specific hyperparameter tuning choices on TRACE's behavior.

3.4 Implications & future directions

Our findings have potentially exciting implications for the discovery of both low-dimensional representations and representations in the original (and anatomically- and/or functionally-relevant, in the case of fMRI) input space. For example, if a study's goal is to induce canonical target patterns of neural activity for a particular object category with real-time decoded neurofeedback (DecNef^{27,28,37}), one might wish to instead 'de-noise' the data by maximizing reconstruction classifier accuracy instead of reconstruction class specificity. In the fMRI dataset presented here, reconstruction classifier accuracy peaked at about $d=200$. It is possible that in other fMRI datasets, reconstruction classifier accuracy might peak at a non-maximal bottleneck dimensionality, in which case it could be used to select the best dimensionality for the task at hand. Alternatively, one could choose to select optimal bottleneck dimensionality based on when reconstruction class specificity or classifier accuracy exceeds the analogous metric calculated directly from the raw input data. Here we showed that TRACE either exceeds these benchmarks sooner than AE, or does so even when AE does not. Thus, the process for selecting the best bottleneck dimensionality can flexibly adapt to an experimenter's goals, and future research seeking to use TRACE to extract neural patterns for use with DecNef should explore how different bottleneck dimensionalities impact the success of the neurofeedback process.

Regardless of the method one uses to select bottleneck dimensionality, though, it seems likely that TRACE can remove task-irrelevant information in a way that is useful for DecNef. To demonstrate this possibility, we did one final exploratory test. Recall that the fMRI dataset used in this study is in part overlapping with the dataset used by Taschereau-Dumouchel and colleagues³⁷, and as such we can directly compare their binary ("cat" versus "everything that is not a cat") decoding accuracy with the decoding accuracy we achieved on TRACE's reconstructions. To translate the reconstruction classifier accuracy we achieved to a binary scale, we counted a prediction to be correct if the correct class was in the top 20 (out of 40) of

predicted classes from our one-versus-all classifier (with chance classification accuracy at 2.5%). Taschereau-Dumouchel and colleagues³⁷ observed binary logistic regression classification accuracies of 71.7% on average within-subject (~1 hour of fMRI data per person). (Relying on hyperalignment⁵³ to pool their 30 subjects and subsequently train such classifiers, they observed mean 82.4% using a 30-subject concatenated dataset). When we trained logistic regression classifiers on each individual subject (i.e., no hyperalignment) – some of whom are actually the original subjects in that former study – and translated the classification accuracies as described to be on the same scale as binary classification, we achieved the equivalent of 94.4% binary accuracy at bottleneck dimensionality $d=30$ (where reconstruction class specificity was maximized). Thus, TRACE facilitates distillation of class-specific representations in native space that are superior to the original representations themselves for this purpose.

Finally, we want to end with a brief note about TRACE’s promise beyond fMRI and DecNef, at a broader scale including other types of biological-scale datasets relevant to human behavior, because fMRI data are not unique in their sparsity of sample size. For example, in biological image analysis or even human microbiome research, we are interested in learning generalizable and biologically informative “truths” about biological systems. However, precisely in the same way as a typical fMRI dataset, biological datasets are often limited in sample size. Given TRACE’s success here, we hope that its capacity to discover task-relevant information *despite* undesirable ratios of samples to input-dimensions can help discover truths about other biological processes. Future studies should apply TRACE to other biological-scale datasets, with the goal of discovering representations relevant to those researchers and domains.

4. Methods

4.1 Methods overview

We proposed the “Task-Relevant Autoencoder via Classifier Enhancement” (TRACE) model and directly compared its behavior to that of a standard autoencoder (AE) with equivalent internal architecture. TRACE is equivalent to the AE model in every respect, with the exception of a classifier branch which reads out directly from the AE network’s bottleneck layer and which contributes to the overall loss function of the network. We benchmarked TRACE and AE on the MNIST and Fashion MNIST datasets often used in machine learning model evaluation, and then applied both models to a previously-collected fMRI dataset to showcase TRACE’s utility in extracting task-relevant low dimensional representations in a nonparametric regime.

We defined four output metrics to quantitatively evaluate the behavior of the TRACE and AE models, described below (**Methods Section 4.4**). These output metrics were evaluated as a function of bottleneck dimensionality -- essentially, how much each dataset can be compressed while retaining task-relevant information in the compressed representation (the bottleneck layer of TRACE and AE). Details of datasets and model architectures are described in the next sections.

4.2 Datasets

4.2.1 MNIST dataset

The MNIST dataset consists of 60,000 handwritten 28x28 pixel grayscale images (i.e., image dimensionality of 784 pixels) of the 10 classes of digits (i.e. 0,1,...9) with their corresponding labels as the training set and 10,000 samples as the test set with a total of 7000 images per

class (training and test set combined) that was collected by LeCun and colleagues²⁵, and is one of the most commonly used benchmarks to evaluate the performance of deep learning models.

4.2.2 Fashion MNIST dataset

The Fashion MNIST dataset²⁶ consists of 60,000 samples of 28x28 grayscale images (i.e., the dimensionality of 784 pixels) of 10 clothing categories (i.e. shirt, shoes, etc) with a total of 7000 images per class (training and test set combined)²⁶.

4.2.3 fMRI dataset

4.2.3.1 Participants & task

The fMRI dataset used here was collected previously for several separate projects and was partially reported previously by Taschereau-Dumouchel and colleagues³⁷. The dataset we used here contained 60 usable subjects' whole-brain data; one subject's data produced wildly unstable outcomes across all metrics tested here (see **Methods Section 4.4**) and so was excluded from the final analysis, and an additional 10 subjects had previously been collected with a different imaging sequence and are not included here. From this dataset, we therefore examined 59 healthy human participants who had viewed 3600 images from 40 categories (90 exemplars of each category, including 30 categories of animals [dogs, cats, snakes, etc.] and 10 categories of man-made objects [keys, chairs, airplanes, etc.]) while whole-brain BOLD responses were acquired. The images were shown for 0.98s each in mini-blocks (chunks) of 2, 3, 4, or 6 images of each category displayed sequentially, and the participants were asked to press a button whenever the category of images was changed. The data was collected from each participant in six runs with short breaks in approximately one hour while they were in the scanner. Each image was on screen for a duration of 0.98 s. The voxel activities of the ventral temporal cortex (VTC) were used as input to the model in order to find the latest features represented by VTC. See the methods reported by Taschereau-Dumouchel and colleagues³⁷ for further details.

4.2.3.2 Image acquisition & preprocessing

The functional imaging data acquisition and preprocessing procedures for this existing dataset are previously described elsewhere³⁷, but are included here for completeness. Participants in the database were scanned at one of two 3T MRI scanners (Siemens Prisma and Verio) with a head coil. Whole brain functional data were acquired in 33 contiguous slices (TR = 2000 ms, TE = 30 ms, voxel size = 3 × 3 × 3.5 mm³, field-of-view = 192 × 192 mm, matrix size = 64 × 64, slice thickness = 3.5 mm, 0 mm slice gap, flip angle = 80 deg) oriented parallel to the AC-PC plane. High-resolution T1-weighted structural MR images (MP-RAGE sequence; 256 slices, TR = 2250 ms, TE = 3.06 ms, 5 voxel size = 1 × 1 × 1 mm³, field-of-view = 256 × 256 mm, matrix size = 256 × 256, slice thickness = 1 mm, 0 mm slice gap, TI = 900 ms, flip angle = 9 deg.) were also obtained. Functional images were preprocessed using standard procedures, including slice timing correction, motion correction, realignment to the first functional image, and coregistration to the structural scan using SPM 12 (Statistical Parametric Mapping; www.fil.ion.ucl.ac.uk/spm). The anatomical mask of the target region of interest, VTC (fusiform, lingual/parahippocampal, and inferior temporal cortex), was selected using the Freesurfer (<http://surfer.nmr.mgh.harvard.edu/>) automated gray matter segmentation combining the ROI labels of *fusiform*, *inferior temporal*, *lingual*, and *parahippocampal*.

Voxels in this combined ventral temporal ROI were detrended and then deconvolved using the least-square separate approach^{57,58}. This method creates an iterative general linear model for

each trial individually, such that the design matrix contains one parameter modeling the current trial, and two parameters modeling all other trials in the design (e.g. even- and odd-numbered trials). This standardized method allows deconvolution of each trial in a rapid-event related design such as this one to obtain parameter estimates for each individual trial. This process results in a $N_{VTC,S} \times 3600$ images (90 exemplars of each category) timeseries, where $N_{VTC,S}$ refers to the number of voxels in the combined VTC ROI for each subject S ; in this dataset $N_{VTC,S}$ ranged from ~40 to 55 MB in CSV format. This timeseries formed the dataset used in this project.

4.3 Models

4.3.1 Standard autoencoder model (AE)

The standard AE model provides a baseline benchmark of model behavior in all datasets tested. This model is designed to be almost as simple as possible in order to maximally highlight any differences between AE and TRACE behaviors. Thus, the standard AE architecture consists of input and output layers with dimensionality equal to the size of the datasets (784 for MNIST and Fashion MNIST, and 1726-3078 for fMRI), plus a single hidden layer with 1000 units in each of the encoding and decoding sections (**Figure 5**).

4.3.1.1 Inputs

To facilitate comparisons across the three datasets tested and to facilitate faster model convergence, we standardized inputs by scaling their values: MNIST and Fashion MNIST inputs were scaled to take on values between 0 and 1 by dividing all values by 255, and fMRI inputs (parameter estimates from the single-trial deconvolution described above) were standardized by z-scoring because these parameter estimates can take on arbitrary real numbers without bound. To prevent leakage of any task-relevant information from test sets to the training set, the scaling factors were determined only on the training set, and then the same scaling parameters were used to scale the test sets.

4.3.1.2 Activation functions

For the hidden layers, we used the hyperbolic tangent as the activation function in order to discover more complex nonlinear patterns in the data, as this function was reported previously to be more sensitive in capturing detailed and local information to represent the data with lower dimensions⁵⁹. For the bottleneck layer of the network, we selected the linear activation function (i.e., no activation function) because in initial explorations we found that using a linear function resulted in discovering more task-relevant features in comparison to other activation functions, such as hyperbolic tangent, rectified linear unit (ReLU) and sigmoid, since features of the bottleneck layer under the linear function had higher accuracy in decoding the categories (e.g. in the case of MNIST dataset and with 2 dimensions in the bottleneck the accuracy increases by 15 percent; full data not shown). For the final decoding layer, a linear activation function was also chosen because the fMRI data are unitless and take arbitrary numbers and therefore are not confined to be within a specific boundary. Using a linear function allows the reconstructing layer to assign any value that minimizes the difference between output and input of the autoencoder, unlike most other typical activation functions (e.g. sigmoid, Tanh, etc) which usually have a confined output value hence those are not good choices to reconstruct fMRI data.

4.3.1.3 Objective function

Equation 1 shows the objective function for the standard autoencoder which was chosen as the mean square error (MSE):

$$L_R = \frac{1}{m \times n} \sum_{i=1}^m \sum_{j=1}^n (\hat{X}_{ij} - X_{ij})^2 \quad (1)$$

Where X is the input with m samples and n input-dimensions, and \hat{X} is the reconstruction of the input.

To minimize the loss function, the Adam⁶⁰ implementation of stochastic gradient descent (SGD) was used and the learning rate was chosen to be 1e-4 and the batch size was set to 32. To prevent overfitting, the dropout technique was used for regularization with a dropout rate of 0.1.

4.3.2 Task-Relevant Autoencoder via Classifier Enhancement (TRACE)

The Task-Relevant Autoencoder via Classifier Enhancement (TRACE) model is identical to the standard AE model with the exception that a logistic regression classifier was attached to the bottleneck (**Figure 9**). The activation function for this “decoder branch” of the network was the softmax function (also known as Boltzman distribution) which outputs a probability distribution for each class (e.g. Classes of 10 digits for MNIST). The loss function for this branch of the network was chosen to be categorical cross-entropy, i.e.:

$$L_{CE} = \frac{-1}{m} \sum_{i=1}^m \sum_{c=1}^k y_{ci} \log(\hat{y}_{ci}) \quad (2)$$

where k denotes the number of the classes, y is the label of observation, and \hat{y} is the predicted label.

In the TRACE network, the final objective function is the summation of reconstruction loss and the categorical cross-entropy loss function (equations 1 and 2), i.e.:

$$\begin{aligned} L_{TRACE} &= L_R + L_{CE} \\ &= \frac{1}{m \times n} \sum_{i=1}^m \sum_{j=1}^n (\hat{X}_{ij} - X_{ij})^2 - \frac{\alpha}{m} \sum_{i=1}^m \sum_{c=1}^k y_{ci} \log(\hat{y}_{ci}) \end{aligned} \quad (3)$$

where α , sets the weight for the classifier part of the loss function in order to control for its participation in updating the parameters.

4.3.3 Training and test data set size

The MNIST and Fashion MNIST datasets, each using 60,000 training and 10,000 test samples, are popular choices in benchmarking deep learning models. Having a high ratio of training samples for the number of data dimensions in these datasets (i.e. $60,000/748 = 76.5$) makes them good candidates even for very deep neural networks. However, verifying the advantage of a model on these huge dataset is not necessarily applicable to dataset with much smaller sample to input-dimension size ratio (e.g. the real-world fMRI dataset we used here with approximate samples-to-input-dimensions ratio of ~ 1.5), as such a ratio might lead to overfitting

and thus poor predictive capacity. To ensure that TRACE is powerful not only when the samples-to-input-dimensions ratio is large but also when available data is much sparser (such as in biological datasets that often suffer from small sample size), we truncated the MNIST and fashion MNIST datasets to explore how TRACE behaves under increasing data truncation.

To accomplish this, we trained both the AE and TRACE models at 10, 30, 50, 70, 90, 95, and 98 percent of data truncation at the optimal bottleneck dimensionality (i.e. $d=2$) for both MNIST and Fashion MNIST. Reducing the number of training samples to 95 and 98 percent decreases the samples-to-input-dimensions ratio such that it is approximately the same as the fMRI dataset also used here. The training set at each level of data truncation was the same for both AE and TRACE networks. We then evaluated the behavior of all four outcome metrics under all levels of data truncation by using the conventional hold-out test set in the MNIST and Fashion MNIST (i.e. 10,000 sample test set) (**see Section 2.1.3**). In the case of the fMRI dataset, we performed the training on 2700 training samples and tested the trained models on 900 held-out test samples. At the extreme level of 98% data truncation in the MNIST and Fashion MNIST datasets, the sample size is reduced to 1200 exemplars with each exemplar having 784 input-dimensions. Therefore, the ratio of samples to input-dimensions is about 1.5 which is comparable to sample-to-input-dimension ratio of fMRI dataset (i.e. ~ 1.5).

4.3.4 Implementation details for TRACE and AE

We explored several different values for α as a hyperparameter (**Equation 3**), and selected 0.01 for all results shown here because this value gave us the lowest MSE loss while having minimal effect on the categorical cross-entropy loss in comparison to higher values α . The total loss was also the lowest when α was set to 0.01 (see **Supplementary Material S1**). For all datasets and all networks we used the same architecture and hyperparameter values. To implement these networks, we used the Keras functional API⁶¹ with TensorFlow⁶² backend, using available GPUs in Google Colab Pro and 12.7 GB of RAM.

4.4 Outcome metrics

In order to explore what is the best low dimensional feature space that explains within class characteristics while preserving the ability of the network to reconstruct the input, we evaluated four metrics as a function of bottleneck dimensionality (2, 5, 10, 15, 20, 30, 50, 60, 100, 150, 200, 250, 500, 1000, or 1500 nodes) in AE and TRACE for all three datasets (MNIST, Fashion MNIST, and fMRI): (1) *reconstruction fidelity*, (2) *reconstruction classifier accuracy*, (3) *bottleneck classifier accuracy*, and (4) *reconstruction class specificity*, as described below and shown in cartoons in **Figure 7**. Collectively, these metrics provide a broad picture of how bottleneck dimensionality can balance reconstruction fidelity with class fidelity, i.e. balances exemplar reconstruction capacity on each trial and per image class (the ‘gist’ of a category) with extraction of task-relevant features in the compressed representation in the bottleneck.

Importantly, the metrics (2) and (5) can be benchmarked against the equivalent metric applied to the *input* (i.e., original data). Specifically, reconstruction classifier accuracy and reconstruction class specificity can be compared to the equivalent assessments for the original, raw input data as a benchmark of the TRACE and AE models’ ability to extract meaningfully less noisy representations and push them back out into the original input dimensionality, such that comparing these metrics between input and reconstruction serves as a meaningful measure of noise reduction accomplished via AE or TRACE.

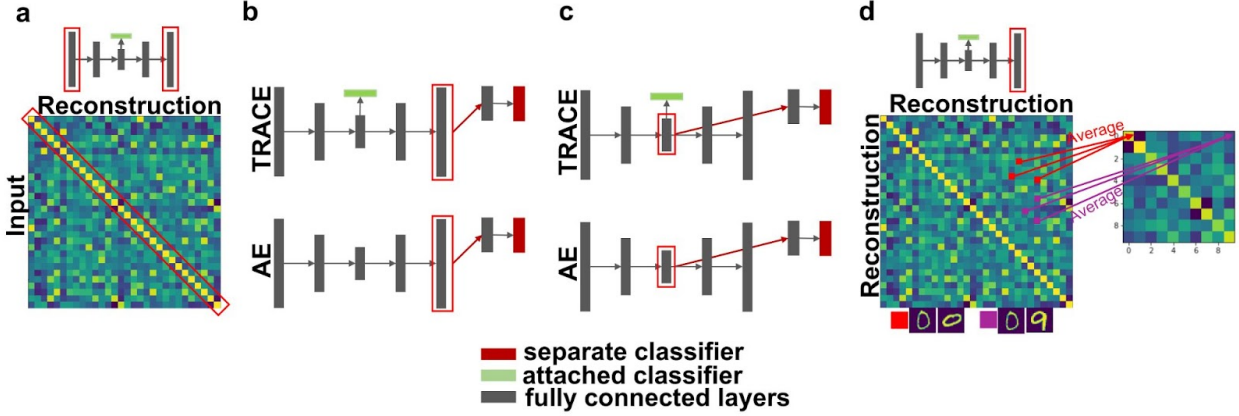


Figure 7. Graphical representation of the four quantitative outcome metrics. **(a)** Reconstruction fidelity (correlation between inputs and reconstructions; **Section 4.4.1**), **(b)** reconstruction classifier accuracy (accuracy of a classifier trained on reconstructions; **Section 4.4.2**), **(c)** bottleneck classifier accuracy (accuracy of a classifier trained on bottleneck features; **Section 4.4.3**), and **(d)** reconstruction class specificity (within- versus between-class separation; **Section 4.4.4**). Small cartoons of the TRACE architecture use red rectangle overlays to indicate which sections of the model architecture are being utilized for each outcome metric. In **(b)** and **(c)**, red-filled boxes indicate separate classifiers (i.e., trained separately after the main network has been trained), green-filled boxes indicate attached classifiers (i.e., which contribute to the main loss function for the network), and gray-filled boxes indicate fully-connected encoder and decoder layers.

4.4.1 Reconstruction fidelity

We quantified how well TRACE and AE could reconstruct the input information with the average of all Pearson correlation coefficients between each input trial of the test set and the corresponding reconstruction of that sample (**Figure 7a**). We computed this correlation coefficient at all bottleneck dimensionalities tested. In the Results, below, this metric is referred to as *reconstruction fidelity* or $fidelity_R$.

$$fidelity_R = E(\rho_R) \quad (4)$$

Where ρ_R is the correlation between each input exemplar and its reconstruction, and E denotes the expected value.

4.4.2 Reconstruction classifier accuracy

To quantify how well the reconstructed input represents a certain class, we used a separate logistic classifier (**Equation 5**) and trained it using reconstructed inputs for all dimensions in the bottleneck (**Figure 7b**). Using the same train/test folds as for training TRACE and AE, we trained the data for 30 epochs for MNIST and Fashion MNIST and 300 epochs for fMRI.

$$L_{RCA} = \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^k y_{ci} \log(\hat{y}_{ci}) + \lambda \sum_{r=1}^p w_r^2 \quad (5)$$

Where L_{RCA} is the cross entropy loss for reconstructed input, and λ is the regularization parameter and w and p are the weight matrices and the number of parameters of the classifier respectively.

4.4.3 Bottleneck classifier accuracy

We quantified the task-relevance of the features in the bottleneck via the accuracy of the logistic regression classifier with such bottleneck node activity as inputs (**Figure 7c**). For both TRACE and AE, this classifier is trained separately, *after* the fully connected network (AE) and fully connected network plus classifier (TRACE) are optimized. That is, this training occurs separately from the global loss function (which also includes an identical logistic regression classifier term, **Equations 2 & 3**) even for TRACE.

To do this, we first extracted the bottleneck features after the training of the TRACE and AE networks was complete, and then trained a completely separate logistic regression classifier to classify the category of the input image as it was now represented in the low-dimensional bottleneck of each model. The loss function for this separate classifier was identical to the one contributing to TRACE's global loss function (categorical cross-entropy), and the learning rate was set to 0.5e-4. To prevent overfitting of the classifier we used the L2 regularization technique. The final objective function for the separate logistic classifier is:

$$L_{BCA} = \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^k y_{ci} \log(\hat{y}_{ci}) + \lambda \sum_{b=1}^q w_b^2 \quad (6)$$

where w and q are the weight matrices and the number of parameters of the classifier respectively. The hyperparameter λ was set to 0.007 which was manually tuned to maximize the classification accuracy.

4.4.4 Reconstruction class specificity

Another measure of the task-relevancy of the reconstructed information is the degree of similarity of representations within a class versus between classes. To this end, we computed the average within-class Pearson correlation of exemplars within each category from the reconstructions by extracting the pairwise correlation of all pairs of within-class exemplars for each category and finding the average of the correlation of them (we excluded the self-correlation of each trial in this calculation since it will always be 1). We then computed the average pairwise Pearson correlation across all pairs of trials from different classes. Graphically, this can be thought of as building an Nclass x Nclass similarity matrix (10x10 for the MNIST and Fashion MNIST data, 40x40 for the fMRI data), in which the diagonal represents the average pairwise within-class similarity (again, excluding self-similarity for a given exemplar) and the off-diagonals (**Figure 7d**). The final measure of representation class specificity from reconstructed inputs is thus the average of the diagonal (within class) of this similarity matrix minus the average of the off-diagonal (between class) of this matrix, i.e.

$$RCS = E(\rho_{R,within}) - E(\rho_{R,between}) \quad (7)$$

where RCS is the class specificity in the reconstruction of the input, and $\rho_{R,within}$ and $\rho_{R,between}$ are the Pearson correlation matrices between trials within each class and between one class

and all other classes based on the reconstructed inputs excluding the elements of the main diagonal of the correlation matrix (i.e. the correlation of trial by itself was excluded since it is always 1 and does not reflect a true with class correlation), and $\rho_{R,between}$ is the Pearson correlation matrix of trials of the reconstructed inputs in one class and trials in all other classes.

We explored how this measure of class specificity changes as the dimension of the bottleneck increases. The purpose of this quantitative analysis is to discover the optimal bottleneck dimensionality for each dataset that can best distinguish within class versus between class categories in the reconstructions (i.e., in the original input space). In **Results** this metric is referred to as *reconstruction class specificity* or *RCS*.

4.4.5 Benchmarks against original inputs

To quantify the reduction in noise and the success of task-relevant feature extraction, we benchmark the reconstructions from AE and TRACE in two ways.

First, we examined the classification accuracy of a simple logistic regression classifier applied to the input data in comparison to the accuracy of an identical classifier applied to the reconstructions (**Methods Section 4.4.2**). That is, if a representation has been successfully de-noised through the compression (and task-relevant feature extraction, in the case of TRACE), then the reduction in task-irrelevant noise should be apparent in the superior classification accuracy of a logistic regression classifier. Thus, we train logistic regression classifiers on the input space as well as the reconstruction (**Methods Section 4.4.2**) at each bottleneck dimensionality, and report the maximum accuracy achieved (Equation 5).

Second, a final test of the ability of TRACE and AE to extract task-relevant representations can be quantified via comparing the reconstruction class specificity (**Methods Section 4.4.4**) against input class specificity, calculated equivalently to reconstruction class specificity (Equation 7).

Data/code availability statement

The data for this project are available from the corresponding authors upon reasonable request. The code implementing the TRACE and AE models, including outcome metrics, is available at [GITHUB LINK AVAILABLE ON ACCEPTANCE].

Author contributions

Syedmehdi Orouji: Conceptualization, analysis, methodology, project administration, software, validation, visualization, writing -- original draft, writing -- review & editing. **Vincent Taschereau-Dumouchel, Aurelio Cortese, Brian Odegaard, Cody Cushing, & Mouslim Cherkaoui:** Data acquisition and preprocessing, validation, visualization, writing -- review & editing. **Mitsuo Kawato & Hakwan Lau:** Conceptualization, validation, visualization, funding acquisition, writing -- review & editing. **Megan A. K. Peters:** Conceptualization, funding acquisition, methodology, project administration, resources, supervision, validation, visualization, writing -- original draft, writing -- review & editing.

Declaration of competing interests

None

Acknowledgements

This work was supported in part by the Canadian Institute for Advanced Research Azrieli Global Scholars Program (M.A.K.P.), the *Fonds de Recherche du Québec - Santé* (V.T-C.), the Innovative Science and Technology Initiative for Security --ATLA (Grant Number JPJ004596; A.c. and M.K.), and JST ERATO (grant number JPMJER1801; A.C. and M.K). Funding sources had no involvement in the design and methodology of the study.

References

1. Sidhu, G. S., Asgarian, N., Greiner, R. & Brown, M. R. G. Kernel Principal Component Analysis for dimensionality reduction in fMRI-based diagnosis of ADHD. *Front. Syst. Neurosci.* **6**, 74 (2012).
2. Mannfolk, P., Wirestam, R., Nilsson, M., Ståhlberg, F. & Olsrud, J. Dimensionality reduction of fMRI time series data using locally linear embedding. *MAGMA* **23**, 327–338 (2010).
3. Yang, Z., LaConte, S., Weng, X. & Hu, X. Ranking and averaging independent component analysis by reproducibility (RAICAR). *Hum. Brain Mapp.* **29**, 711–725 (2008).
4. Chen, P.-H. (cameron) *et al.* A Reduced-Dimension fMRI Shared Response Model. in *Advances in Neural Information Processing Systems* (eds. Cortes, C., Lawrence, N., Lee, D., Sugiyama, M. & Garnett, R.) vol. 28 460–468 (Curran Associates, Inc., 2015).
5. Haxby, J. V. *et al.* Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science* **293**, 2425–2430 (2001).
6. Haynes, J.-D. & Rees, G. Decoding mental states from brain activity in humans. *Nat. Rev. Neurosci.* **7**, 523–534 (2006).
7. Norman, K. A., Polyn, S. M., Detre, G. J. & Haxby, J. V. Beyond mind-reading: multi-voxel pattern analysis of fMRI data. *Trends Cogn. Sci.* **10**, 424–430 (2006).
8. Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 8619–8624 (2014).
9. Chen, X. *et al.* InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. in *Advances in Neural Information Processing Systems* (eds. Lee, D., Sugiyama, M., Luxburg, U., Guyon, I. & Garnett, R.) vol. 29 (Curran Associates, Inc., 2016).
10. Schölkopf, B., Smola, A. & Müller, K.-R. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural Comput.* **10**, 1299–1319 (1998).

11. Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**, 559–572 (1901).
12. Eckart, C. & Young, G. The approximation of one matrix by another of lower rank. *Psychometrika* **1**, 211–218 (1936).
13. Bejjanki, V. R., da Silveira, R. A., Cohen, J. D. & Turk-Browne, N. B. Noise correlations in the human brain and their impact on pattern classification. *PLoS Comput. Biol.* **13**, e1005674 (2017).
14. Liu, T. T. Noise contributions to the fMRI signal: An overview. *Neuroimage* **143**, 141–151 (2016).
15. Peltier, S. J. *Characterization and Compensation of Systematic Noise in Functional Magnetic Resonance Imaging*. (University of Michigan., 2003).
16. Braun, M. L., Buhmann, J. M. & Müller, K.-R. On relevant dimensions in kernel feature spaces. *J. Mach. Learn. Res.* **9**, 1875–1908 (2008).
17. Cox, D. D. & Savoy, R. L. Functional magnetic resonance imaging (fMRI) ‘brain reading’: detecting and classifying distributed patterns of fMRI activity in human visual cortex. *NeuroImage* vol. 19 261–270 (2003).
18. Wasikowski, M. & Chen, X.-W. Combating the Small Sample Class Imbalance Problem Using Feature Selection. *IEEE Trans. Knowl. Data Eng.* **22**, 1388–1400 (2010).
19. He, H. & Garcia, E. A. Learning from Imbalanced Data. *IEEE Trans. Knowl. Data Eng.* **21**, 1263–1284 (2009).
20. Nie, D., Zhang, H., Adeli, E., Liu, L. & Shen, D. 3D Deep Learning for Multi-modal Imaging-Guided Survival Time Prediction of Brain Tumor Patients. *Med. Image Comput. Comput. Assist. Interv.* **9901**, 212–220 (2016).
21. Socher, R., Pennington, J., Huang, E. H., Ng, A. Y. & Manning, C. D. Semi-Supervised Recursive Autoencoders for Predicting Sentiment Distributions. in *Proceedings of the 2011*

- Conference on Empirical Methods in Natural Language Processing* 151–161 (Association for Computational Linguistics, 2011).
22. Ghifary, M., Kleijn, W. B., Zhang, M., Balduzzi, D. & Li, W. Deep Reconstruction-Classification Networks for Unsupervised Domain Adaptation. in *Computer Vision – ECCV 2016* 597–613 (Springer International Publishing, 2016).
 23. Li, X. *et al.* Intelligent cross-machine fault diagnosis approach with deep auto-encoder and domain adaptation. *Neurocomputing* **383**, 235–247 (2020).
 24. Hosoya, H. CIGMO: Learning categorical invariant deep generative models from grouped data. (2020).
 25. LeCun, Y., Cortes, C. & Burges, C. MNIST handwritten digit database, 1998. URL <http://www.research.att.com/~yann/ocr/mnist> (1998).
 26. Xiao, H., Rasul, K. & Vollgraf, R. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *arXiv [cs.LG]* (2017).
 27. Shibata, K., Sasaki, Y., Kawato, M. & Watanabe, T. Perceptual learning incepted by decoded fMRI neurofeedback without stimulus presentation. *Journal of Vision* vol. 12 282–282 (2012).
 28. Watanabe, T., Sasaki, Y., Shibata, K. & Kawato, M. Advances in fMRI Real-Time Neurofeedback. *Trends Cogn. Sci.* **21**, 997–1010 (2017).
 29. Amano, K., Shibata, K., Kawato, M., Sasaki, Y. & Watanabe, T. Learning to Associate Orientation with Color in Early Visual Areas by Associative Decoded fMRI Neurofeedback. *Curr. Biol.* **26**, 1861–1866 (2016).
 30. Koizumi, A. *et al.* Fear reduction without fear through reinforcement of neural activity that bypasses conscious exposure. *Nature Human Behaviour* **1**, 1–7 (2016).
 31. Shibata, K. *et al.* Toward a comprehensive understanding of the neural mechanisms of decoded neurofeedback. *Neuroimage* **188**, 539–556 (2019).
 32. LaConte, S. M., Peltier, S. J. & Hu, X. P. Real-time fMRI using brain-state classification.

- Hum. Brain Mapp.* **28**, 1033–1044 (2007).
33. deCharms, R. C. *et al.* Learned regulation of spatially localized brain activation using real-time fMRI. *Neuroimage* **21**, 436–443 (2004).
 34. Scharnowski, F., Hutton, C., Josephs, O., Weiskopf, N. & Rees, G. Improving visual perception through neurofeedback. *J. Neurosci.* **32**, 17830–17841 (2012).
 35. Cortese, A., Amano, K., Koizumi, A., Kawato, M. & Lau, H. Multivoxel neurofeedback selectively modulates confidence without changing perceptual performance. *Nat. Commun.* **7**, 13669 (2016).
 36. Scheinost, D. *et al.* Orbitofrontal cortex neurofeedback produces lasting changes in contamination anxiety and resting-state connectivity. *Transl. Psychiatry* **3**, e250 (2013).
 37. Taschereau-Dumouchel, V. *et al.* Towards an unconscious neural reinforcement intervention for common fears. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 3470–3475 (2018).
 38. Shibata, K., Watanabe, T., Kawato, M. & Sasaki, Y. Differential Activation Patterns in the Same Brain Region Led to Opposite Emotional States. *PLoS Biol.* **14**, e1002546 (2016).
 39. Cortese, A., Amano, K., Koizumi, A., Lau, H. & Kawato, M. Decoded fMRI neurofeedback can induce bidirectional confidence changes within single participants. *NeuroImage* vol. 149 323–337 (2017).
 40. Caballero-Gaudes, C. & Reynolds, R. C. Methods for cleaning the BOLD fMRI signal. *Neuroimage* **154**, 128–149 (2017).
 41. Shmueli, K. *et al.* Low-frequency fluctuations in the cardiac rate as a source of variance in the resting-state fMRI BOLD signal. *Neuroimage* **38**, 306–320 (2007).
 42. Birn, R. M., Diamond, J. B., Smith, M. A. & Bandettini, P. A. Separating respiratory-variation-related fluctuations from neuronal-activity-related fluctuations in fMRI. *Neuroimage* **31**, 1536–1548 (2006).
 43. Fullana, M. A. *et al.* Neural signatures of human fear conditioning: an updated and extended meta-analysis of fMRI studies. *Molecular Psychiatry* vol. 21 500–508 (2016).

44. Hofmann, S. G., Ellard, K. K. & Siegle, G. J. Neurobiological correlates of cognitions in fear and anxiety: A cognitive–neurobiological information-processing model. *Cognition and Emotion* **26**, 282–299 (2012).
45. Hinton, G. E. & Salakhutdinov, R. R. Reducing the dimensionality of data with neural networks. *Science* **313**, 504–507 (2006).
46. Wang, S., Ding, Z. & Fu, Y. Coupled Marginalized Auto-Encoders for Cross-Domain Multi-View Learning. in *IJCAI* 2125–2131 (2016).
47. Kingma, D. P., Mohamed, S., Rezende, D. J. & Welling, M. Semi-supervised learning with deep generative models. in *Advances in neural information processing systems* 3581–3589 (2014).
48. Maaløe, L., Sønderby, C. K., Sønderby, S. K. & Winther, O. Improving semi-supervised learning with auxiliary deep generative models. in *NIPS Workshop on Advances in Approximate Bayesian Inference* (2015).
49. Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I. & Frey, B. Adversarial Autoencoders. *arXiv [cs.LG]* (2015).
50. Goodfellow, I. *et al.* Generative Adversarial Nets. in *Advances in Neural Information Processing Systems* (eds. Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. & Weinberger, K. Q.) vol. 27 (Curran Associates, Inc., 2014).
51. Radford, A., Metz, L. & Chintala, S. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *arXiv [cs.LG]* (2015).
52. Wang, S., Ding, Z. & Fu, Y. Feature Selection Guided Auto-Encoder. *AAAI* **31**, (2017).
53. Haxby, J. V. *et al.* A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron* **72**, 404–416 (2011).
54. Guntupalli, J. S. *et al.* A Model of Representational Spaces in Human Cortex. *Cereb. Cortex* bhw068–bhw068 (2016).
55. Busch, E. L. *et al.* Hybrid Hyperalignment: A single high-dimensional model of shared

- information embedded in cortical patterns of response and functional connectivity. *Cold Spring Harbor Laboratory* 2020.11.25.398883 (2020) doi:10.1101/2020.11.25.398883.
56. Huang, J. *et al.* Learning shared neural manifolds from multi-subject fMRI data. *arXiv [q-bio.NC]* (2021).
 57. Mumford, J. A., Turner, B. O., Ashby, F. G. & Poldrack, R. A. Deconvolving BOLD activation in event-related designs for multivoxel pattern classification analyses. *Neuroimage* **59**, 2636–2643 (2012).
 58. Turner, B. O., Mumford, J. A., Poldrack, R. A. & Ashby, F. G. Spatiotemporal activity estimation for multivoxel pattern analysis with rapid event-related designs. *Neuroimage* **62**, 1429–1438 (2012).
 59. Ng, W. W. Y., Zeng, G., Zhang, J., Yeung, D. S. & Pedrycz, W. Dual autoencoders features for imbalance classification problem. *Pattern Recognit.* **60**, 875–889 (2016).
 60. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv [cs.LG]* (2014).
 61. Chollet, F. & Others. Keras documentation. *keras. io* **33**, (2015).
 62. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *arXiv [cs.DC]* (2016).

Supplementary Material

S1 Hyperparameter tuning

In order to tune the hyperparameter α to control for the contribution of the logistic classifier to the final objective function of TRACE, we chose different values for α (i.e. 0, 0.01, 0.1, 0.2, 0.5, 0.9, 1) for the bottleneck dimensionality of $d=2$. We chose α as 0.01 since it seemed it is the optimum point for the reconstruction and classification trade-off. Surprisingly, at $\alpha=0.01$ reconstruction loss was actually lower than $\alpha=0$ and the total loss was found to be almost the same as $\alpha=0$.

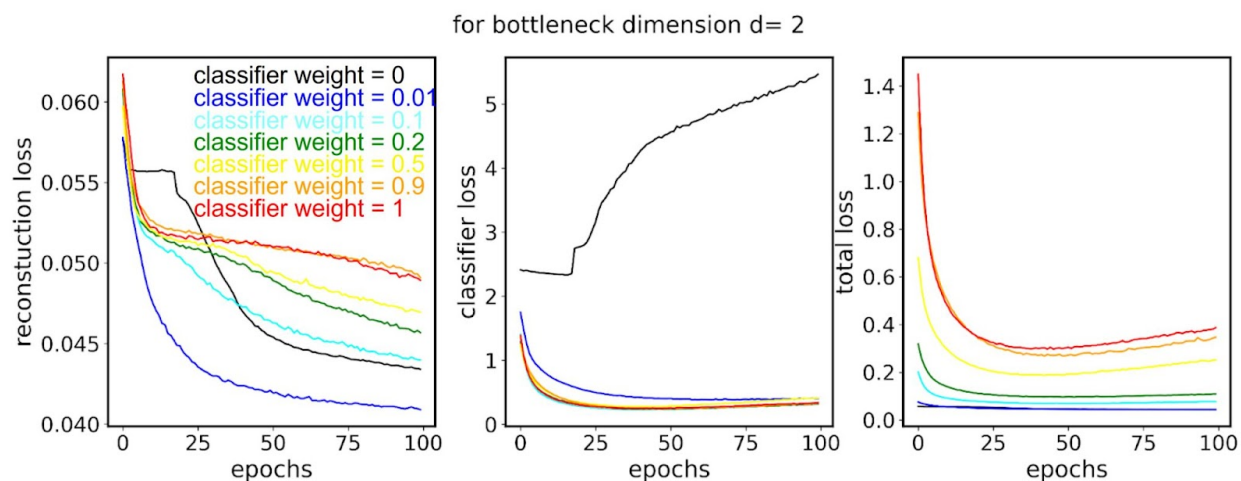


Figure S1. Loss functions for the MNIST dataset as a function of training epochs for bottleneck dimensionality $d=2$, for different levels of the hyperparameter α . We selected $\alpha=0.01$ for all results presented in the main text because it showed the fastest convergence and achieved the smallest total loss (bottom panel) when training was complete.

S2 Model fitting practical details

In order to meaningfully compare behaviors between the TRACE and AE models, it is important to determine that both models can both be adequately fit to each of our datasets. Using available GPU processors in Google Colab Pro, it took about 61 minutes to fit the AE model to the MNIST training dataset (60,000 labeled samples), and about the same time for Fashion MNIST training dataset (60,000 labeled samples), and 21 minutes on average for each human subject in the fMRI dataset (3600 labeled samples of VTC) and for all 15 dimensions of bottleneck.

S3 Comprehensive discussion of model performance

In the main text, we present and discuss the results from bottleneck dimensionalities between 2 - 250, because at $d > 250$ the quantitative metrics (see Methods) tend to asymptote. **Figure S2** presents the four metrics at bottleneck dimensionalities up to 1500 as comprehensive demonstration of this behavior. **Table S1** provides the bottleneck dimensionality at which each metric is maximized.

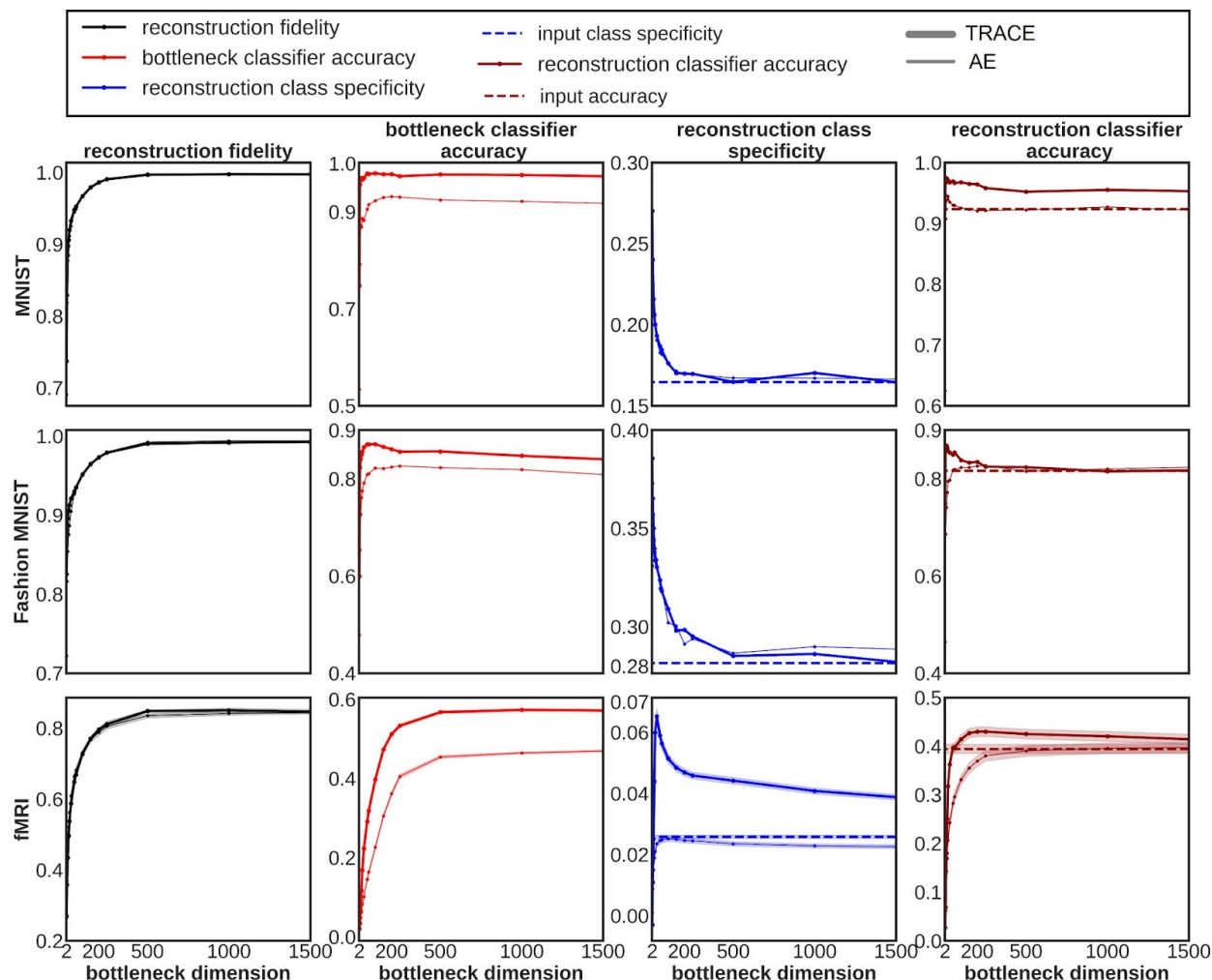


Figure S2. Quantitative comparison between TRACE and AE on the four outcome metrics, for all three datasets (MNIST, Fashion MNIST, and fMRI) for bottleneck dimensionalities between 2 and 1500. As in Figure 3, the black, red, and blue TRACES show the reconstruction fidelity ($fidelity_R$), bottleneck classifier accuracy (BCA), and reconstruction class specificity (RCS), respectively. Lighter TRACES in the fMRI dataset results (bottom row) represent individual subjects, with the thicker TRACE showing the mean of each metric across subjects. The red dashed rectangles show the portions of the plots shown in the main text.

			MNIST	Fashion MNIST	fMRI
1	Reconstruction fidelity	AE	(no peak)	(no peak)	(no peak)
		TRACE	(no peak)	(no peak)	(no peak)
2	Reconstruction classifier accuracy	AE	20	(no peak)	(no peak)
		TRACE	10	10	asymptote at ~150
3	Bottleneck classifier accuracy	AE	200	200	500
		TRACE	15	50	500
4	Reconstruction class specificity	AE	10	5	100
		TRACE	2	2	30

Table S1. Bottleneck dimensionality at which each outcome metric is maximized.

			MNIST	Fashion MNIST	fMRI
1	Reconstruction fidelity	AE	0.63	0.78	0.61 ± 0.05
		TRACE	0.67	0.79	0.59 ± 0.08
2	Reconstruction classifier accuracy	AE	0.47	0.60	0.24 ± 0.07
		TRACE	0.77	0.68	0.36 ± 0.07
3	Bottleneck classifier accuracy	AE	0.39	0.50	0.1 ± 0.04
		TRACE	0.72	0.63	0.22 ± 0.04
4	Reconstruction class specificity	AE	0.11	0.30	0.02 ± 0.006
		TRACE	0.16	0.31	0.07 ± 0.02

Table S2. Performance at bottleneck dimensionality d=2 for MNIST and Fashion MNIST, and d=30 for fMRI.

S4 Principal components analysis

One might wonder whether linear PCA would produce similarly appropriate reconstructions to either TRACE or AE. To test this possibility, we directly compared the reconstruction quality and clustering behavior for TRACE and AE to the output from the top two components of linear PCA. The results demonstrate that the reconstruction and clustering behavior for the top two principal components from linear PCA is woefully inadequate (**Figure S3**). Thus, TRACE is superior to linear PCA by these metrics.

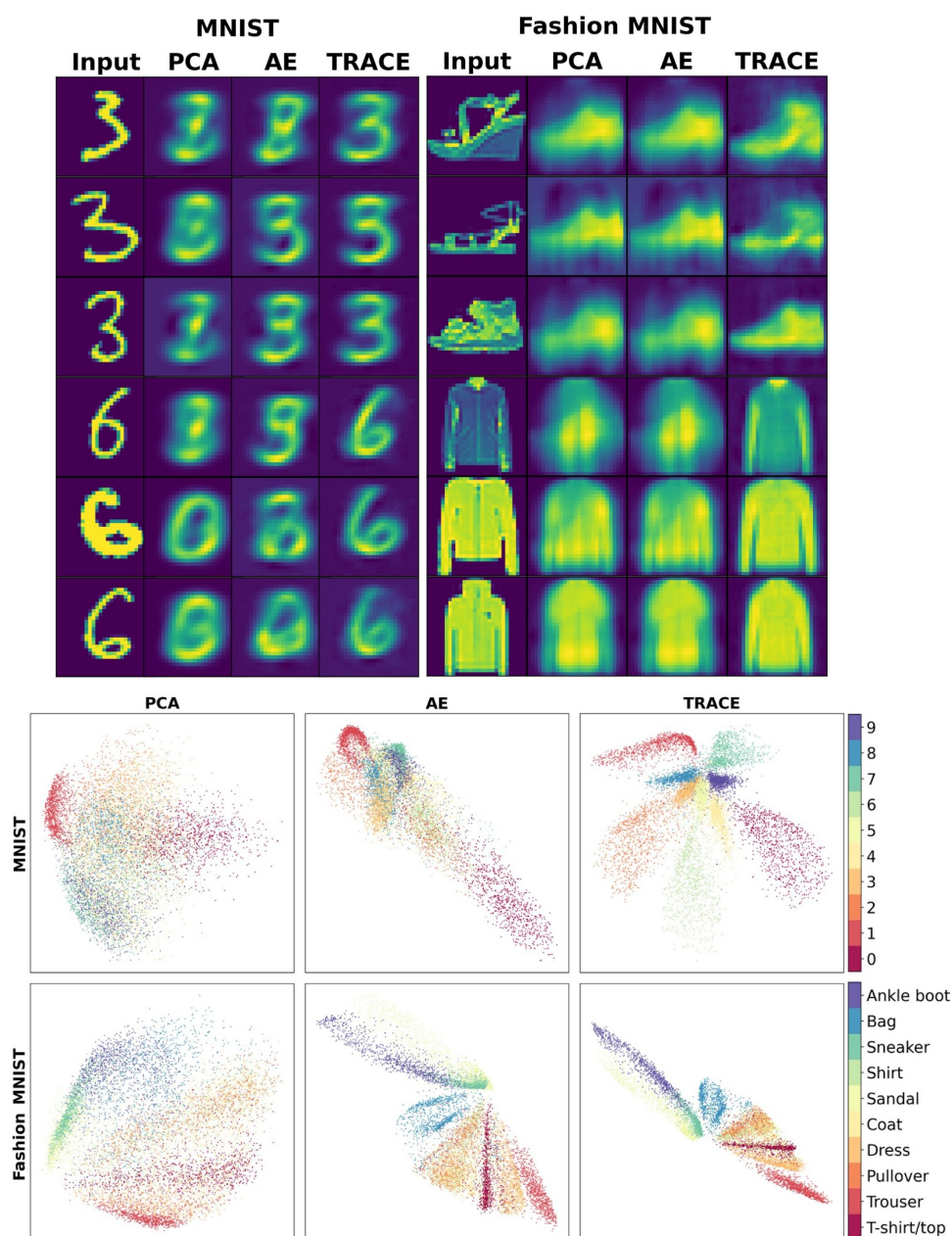


Figure S3. Comparison of AE and TRACE behaviors to linear PCA, for the MNIST and Fashion MNIST datasets.