
Data-Efficient Mutual Information Neural Estimator

Xiao Lin^{1*}, Indranil Sur^{1*}, Samuel A. Nastase²,
Ajay Divakaran¹, Uri Hasson² and Mohamed R. Amer¹

¹SRI International, Princeton, NJ, USA ²Princeton University, Princeton, NJ, USA

Abstract

Measuring Mutual Information (MI) between high-dimensional, continuous, random variables from observed samples has wide theoretical and practical applications. Recent work, *MINE* [5], focused on estimating tight variational lower bounds of MI using neural networks, but assumed unlimited supply of samples to prevent overfitting. In real world applications, data is not always available at a surplus. In this work, we focus on improving data efficiency and propose a Data-Efficient MINE Estimator (DEMINE), by developing a relaxed predictive MI lower bound that can be estimated at higher data efficiency by orders of magnitudes. The predictive MI lower bound also enables us to develop a new meta-learning approach using task augmentation, Meta-DEMINE, to improve generalization of the network and further boost estimation accuracy empirically. With improved data-efficiency, our estimators enables statistical testing of dependency at practical dataset sizes. We demonstrate the effectiveness of our estimators on synthetic benchmarks and a real world fMRI data, with application of inter-subject correlation analysis.

1 Introduction

Mutual Information (MI) is an important, theoretically grounded, measure of similarity between random variables. MI captures general, non-linear, statistical dependencies between random variables. It is a widely used quantity in various machine learning tasks ranging from classification to feature selection and neural network analysis.

A widely used approach for estimating MI from samples is using k-NN estimates, notably the KSG estimator [29]. [15] provided a comprehensive review and studied the consistency and of asymptotic confidence bound of the KSG estimator [16]. MI estimation can also be achieved by estimating individual entropy terms involved through kernel density estimation [2] or cross-entropy [31]. Overfitting can be reduced through partitioning the samples into different folds for modeling and for estimation. Despite of their fast and accurate estimations on random variables with few dimensions, MI estimation on high-dimensional random variables remains challenging for commonly used Gaussian kernels. Fundamentally, estimating MI requires the ability to accurately model the random variables, where high-capacity neural networks have shown excellent performance on complex high-dimensional signals such as text, image and audio.

Recent works on MI estimation have focused on developing tight variational MI lower bounds where neural networks are used for signal modeling. The IM algorithm [1] introduces a variational MI lower bound, where a neural network $q(z|x)$ is learned as a variational approximation to the conditional distribution $P(Z|X)$. The IM algorithm requires the entropy, $H(Z)$, and $E_{XZ} \log q(z|x)$ to be tractable, which applies to latent codes of Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs) as well as categorical variables. [5] introduces MI lower bounds *MINE* and *MINE-f* which allow the modeling of general random variables and shows improved accuracy for high-dimensional random variables, with application to improving generative models. [35] introduces

*equal contribution

a spectrum of energy-based MI estimators based on *MINE* and *MINE-f* lower bounds and a new TCPC estimator for the case when multiple samples from $P(Z|X)$ can be drawn.

An important challenge that previous works overlooked is MI estimation using limited data. As the high-capacity neural networks tend to overfit. Variational estimators, such as *MINE*, expect an impractically large number of samples to overcome overfitting and to reach high confidence. In addition, tighter lower bounds may also require more data to estimate. When limited number of samples are provided, estimations can suffer from high variance observed in [35].

To address the data efficiency challenge, our estimator, DEMINE, introduces predictive mode and meta-learning to the *MINE* estimator family to greatly improve sample efficiency. We develop a relaxed, predictive variational lower bound based on *MINE* that prevents overfitting by explicitly partitioning samples into training and validation. Furthermore, a predictive formulation allows us to incorporate techniques that improves generalization beyond curve fitting such as meta-learning. With these improvements, we show that DEMINE enables practical statistical testing of dependency in not only synthetic datasets but also for real world functional Magnetic Resonance Imaging (fMRI) data analysis for capturing nonlinear and higher-order brain-to-brain coupling.

An additional component to enhance our estimators is meta-learning. Meta-learning, or "learning to learn", seeks to improve the generalization capability of neural networks by searching for better hyper parameters [30], network architectures [34], initialization [10, 11, 27] and distance metrics [46, 42]. Meta-learning approaches have shown significant performance improvements in applications such as automatic neural architecture search [34], few-shot image recognition [10] and imitation learning [12].

In particular, our estimator benefits from the Model-Agnostic Meta-Learning (MAML) [10] framework which is designed to improve few-shot learning performance. A network initialization is learned to maximize its performance when fine-tuned on few-shot learning tasks. Applications include few-shot image classification and navigation. We leverage the model-agnostic nature of MAML for MI estimation between generic random variable and adopt MAML for maximizing MI lower bounds. To construct a collection of diverse tasks for MAML learning from limited samples, inspired by MI's invariance to invertible transformations, we propose a task-augmentation protocol to automatically construct tasks by sampling random transformations to transform the samples. Results show reduced overfitting and improved generalization.

Our contributions are summarized as follows: 1) Data Efficient Mutual Information Neural Estimator (DEMINE); 2) New formulation of meta-learning using Task Augmentation (Meta-DEMINE); 3) Application to real life, data scarce application (fMRI).

2 Background

In this section, we will provide the background necessary to understand our approach². We define X and Z to be two random variables, $P(X, Z)$ is the joint distribution, and $P(X)$ and $P(Z)$ are the marginal distributions over X and Z respectively. Our goal is to estimate MI, $I(X; Z)$ given *i.i.d.* sample pairs (x_i, z_i) , $i = 1, 2 \dots n$ from $P(X, Z)$. Let $\mathcal{F} = \{T_\theta(x, z)\}_{\theta \in \Theta}$ be a class of scalar functions, where θ is the set of model parameters. Let $q(x|z) = p(x) \frac{e^{T_\theta(x, z)}}{\mathbb{E}_{(x, z) \sim P_{XZ}} e^{T_\theta(x, z)}}$. the following energy-based family of lower bounds of MI hold for any θ :

$$\begin{aligned} I(X; Z) &\geq \mathbb{E}_{(x, z) \sim P_{XZ}} \log \frac{q(x|z)}{p(x)} = \mathbb{E}_{(x, z) \sim P_{XZ}} T_\theta(x, z) - \mathbb{E}_{x \sim P_X} \log \mathbb{E}_{z \sim P_Z} e^{T_\theta(x, z)} \triangleq I_{EB1} \quad [35] \\ &\geq \mathbb{E}_{(x, z) \sim P_{XZ}} T_\theta(x, z) - \log \mathbb{E}_{x \sim P_X, z \sim P_Z} e^{T_\theta(x, z)} \triangleq I_{MINE} \quad [5] \\ &\geq \mathbb{E}_{(x, z) \sim P_{XZ}} T_\theta(x, z) - \mathbb{E}_{x \sim P_X, z \sim P_Z} e^{T_\theta(x, z)} + 1 \triangleq I_{MINE-f} \quad [5], I_{EB} \quad [35] \end{aligned} \tag{1}$$

where, \mathbb{E} is the expectation over the given distribution. Based on I_{MINE} , the *MINE* estimator $\widehat{I(X; Z)}_n$ is defined as in Eq.2. Estimators for I_{EB1} , I_{MINE-f} and I_{EB} can be defined similarly.

$$\widehat{I(X; Z)}_n = \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n T_\theta(x_i, z_i) - \log \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n e^{T_\theta(x_i, z_j)}. \tag{2}$$

With infinite samples to approximate expectation, Eq.2 converges to the lower bound $\widehat{I(X; Z)}_\infty = \sup_{\theta \in \Theta} I_{MINE}$. Note that the number of samples n needs to be substantially more than the number of model parameters $d = |\theta|$ to prevent $T_\theta(X, Y)$ from overfitting to the samples (x_i, z_i) , $i = 1, 2 \dots n$

²We follow the same notation in [5]. We encourage the review of [5, 35] to understand I_{MINE} , I_{EB1} , and I_{EB} .

and overestimating MI. Formally, the sample complexity of *MINE* is defined as the minimum number of samples n in order to achieve Eq.3,

$$\Pr(|\widehat{I(X, Z)}_n - \widehat{I(X, Z)}_\infty| \leq \epsilon) \geq 1 - \delta. \quad (3)$$

Specifically, *MINE* proves that under the following assumptions: 1) $T_\theta(X, Z)$ is L -Lipschitz; 2) $T_\theta(X, Z) \in [-M, M]$, 3) $\{\theta_i \in [-K, K], \forall i \in 1, \dots, d\}$, the sample complexity of *MINE* is given by Eq.4.

$$n \geq \frac{2M^2(d \log(16KL\sqrt{d}/\epsilon) + 2dM + \log(2/\delta))}{\epsilon^2}. \quad (4)$$

For example, a neural network with dimension $d = 10,000$, $M = 1$, $K = 0.1$ and $L = 1$, achieving a confidence interval of $\epsilon = 0.1$ with 95% confidence would require $n \geq 18,756,256$ samples. This is achievable for synthetic example generated by Generative Adversarial Networks (GANs). For real data, however, the cost of data acquisition for reaching statistically significant estimation can be prohibitively expensive. We propose to use the MI lower bounds specified in Eq.1 from a prediction perspective, inspired by cross-validation. Our estimator, DEMINE, improves sample complexity by disentangling data for lower bound estimation from data for learning a generalizable $T_\theta(X, Z)$. DEMINE enables high-confidence MI estimation on small datasets. Bound tightness is further improved by Meta-DEMINE by using meta-learning to learn generalizable $T_\theta(X, Z)$.

3 Approach

§3.1 specifies DEMINE for predictive MI estimation and derives the confidence interval; §3.2 formulates Meta-DEMINE, explains task augmentation, and defines the optimization algorithms.

3.1 Predictive Mutual Information Estimation

In DEMINE, we interpret the estimation of *MINE*- f lower bound³ Eq.1 as a learning problem. The goal is to infer the optimal network $T_{\theta^*}(X, Z)$ with parameters θ^* using a limited number of samples defined as follows:

$$\theta^* = \arg \max_{\theta \in \Theta} \mathbb{E}_{P_{XZ}} T_\theta(X, Z) - \mathbb{E}_{P_X} \mathbb{E}_{P_Z} e^{T_\theta(X, Z)} + 1.$$

Specifically, samples from $P(X, Z)$ are subdivided into a training set $\{(x_i, z_i)_{\text{train}}, i = 1, \dots, m\}$ and a validation set $\{(x_i, z_i)_{\text{val}}, i = 1, \dots, n\}$. The training set is used for learning a network $\tilde{\theta}$ as an approximation to θ^* whereas the validation set is used for computing the DEMINE estimation $\widehat{I(X, Z)}_{n, \tilde{\theta}}$ defined as in Eq.5.

$$\widehat{I(X, Z)}_{n, \tilde{\theta}} = \frac{1}{n} \sum_{i=1}^n T_{\tilde{\theta}}(x_i, z_i)_{\text{val}} - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n e^{T_{\tilde{\theta}}(x_i, z_j)_{\text{val}}} + 1 \quad (5)$$

We propose a approach to learn $\tilde{\theta}$, DEMINE. DEMINE learns $\tilde{\theta}$ by maximizing the MI lower bound on the training set as follows:

$$\tilde{\theta} = \arg \min_{\theta \in \Theta} \mathcal{L}(\{(x, z)\}_{\text{train}}, \theta), \text{ where,}$$

$$\mathcal{L}(\{(x, z)\}_{\mathcal{B}}, \theta) = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} T_\theta(x_i, z_i)_{\mathcal{B}} + \frac{1}{|\mathcal{B}|^2} \sum_{i=1}^{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} e^{T_\theta(x_i, z_j)_{\mathcal{B}}} - 1. \quad (6)$$

The DEMINE algorithm is shown in Algorithm 1.

³*MINE* lower bound can also be interpreted in the predictive way, but will result in a higher sample complexity than *MINE*- f lower bound. We choose *MINE*- f in favor of a lower sample complexity over bound tightness.

Algorithm 1 DEMINE

Input Data: $\{(x, z)_{\text{train}}, (x, z)_{\text{val}}\}$
Parameters: Batch \mathcal{B} , Iterations N_O , Learning rate η
Output: MI, $T_{\tilde{\theta}}(X, Z)$
 1: $\theta^{(0)} \leftarrow$ Xavier Initialization [18]
 2: **for** $i = 1 : N_O$ **do**
 3: Sample a batch of $(x_i, z_i)_{\mathcal{B}} \sim (x, z)_{\text{train}}$
 4: Compute $\mathcal{L} \left((x_i, z_i)_{\mathcal{B}}, \theta^{(i-1)} \right)$
 5: Compute $\nabla_{\theta}^{(i)} \mathcal{L}$ – gradient for θ
 6: Update $\theta^{(i)}$ using Adam with η
 7: **end for**
 8: MI = $\widehat{I(X, Z)}_{n, \theta^{(N_O)}}$
 9: **return** MI, $\theta^{(N_O)}$

Sample complexity analysis. Because $\tilde{\theta}$ is learned independently of validation samples $\{(x_i, z_i)_{\text{val}}, i = 1, \dots, n\}$, the sample complexity of the DEMINE estimator does not involve the model class \mathcal{F} and the sample complexity is greatly reduced compared to *MINE-f*. DEMINE estimates $\widehat{I(X, Z)}_{\infty, \tilde{\theta}}$ when infinite number of samples are provided, defined as:

$$\begin{aligned} \widehat{I(X, Z)}_{\infty, \tilde{\theta}} &= \mathbb{E}_{P_{XZ}} T_{\tilde{\theta}}(X, Z) - \mathbb{E}_{P_X} \mathbb{E}_{P_Z} e^{T_{\tilde{\theta}}(X, Z)} + 1 \\ &\leq \sup_{\theta \in \Theta} \mathbb{E}_{P_{XZ}} T_{\theta}(X, Z) - \mathbb{E}_{P_X} \mathbb{E}_{P_Z} e^{T_{\theta}(X, Z)} + 1 \leq I(X; Z) \end{aligned} \quad (7)$$

We now derive the sample complexity of DEMINE defined as the number of samples n required for $\widehat{I(X, Z)}_{n, \tilde{\theta}}$ to be a good approximation to $\widehat{I(X, Z)}_{\infty, \tilde{\theta}}$ in Theorem 1.

Theorem 1. For $T_{\tilde{\theta}}(X, Z)$ bounded by $[L, U]$, given any accuracy ϵ and confidence δ , we have:

$$\Pr(|\widehat{I(X, Z)}_{n, \tilde{\theta}} - \widehat{I(X, Z)}_{\infty, \tilde{\theta}}| \leq \epsilon) \geq 1 - \delta$$

when the number of validation samples n satisfies:

$$n \geq n^*, \text{ s.t. } f(n^*) \equiv \min_{0 \leq \xi \leq \epsilon} 2e^{-\frac{2\xi^2 n^*}{(U-L)^2}} + 4e^{-\frac{(\epsilon-\xi)^2 n^*}{2(\epsilon^U - \epsilon^L)^2}} = \delta \quad (8)$$

Proof. Since $T_{\tilde{\theta}}(X, Z)$ is bounded by $[L, U]$, applying the Hoeffding inequality to the first half of Eq.5:

$$\Pr\left(\left|\frac{1}{n} \sum_{i=1}^n T_{\tilde{\theta}}(x_i, z_i) - \mathbb{E}_{P_{XZ}} T_{\tilde{\theta}}(X, Z)\right| \geq \xi\right) \leq 2e^{-\frac{2\xi^2 n}{(U-L)^2}}$$

As $e^{T_{\tilde{\theta}}(X, Z)}$ is bounded by $[e^L, e^U]$, applying the Hoeffding inequality to the second half of Eq.5:

$$\begin{aligned} \Pr\left(\left|\mathbb{E}_{P_X} \mathbb{E}_{P_Z} e^{T_{\tilde{\theta}}(X, Z)} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{P_Z} e^{T_{\tilde{\theta}}(x_i, z_i)}\right| \geq \zeta\right) &\leq 2e^{-\frac{2\zeta^2 n}{(\epsilon^U - \epsilon^L)^2}} \\ \Pr\left(\left|\mathbb{E}_{P_Z} \frac{1}{n} \sum_{i=1}^n e^{T_{\tilde{\theta}}(x_i, z_i)} - \frac{1}{n} \sum_{j=1}^n \frac{1}{n} \sum_{i=1}^n e^{T_{\tilde{\theta}}(x_i, z_j)}\right| \geq \zeta\right) &\leq 2e^{-\frac{2\zeta^2 n}{(\epsilon^U - \epsilon^L)^2}} \end{aligned}$$

Combining the above bounds results in:

$$\Pr(|\widehat{I(X, Z)}_{n, \tilde{\theta}} - \widehat{I(X, Z)}_{\infty, \tilde{\theta}}| \leq \xi + 2\zeta) \geq 1 - 2e^{-\frac{2\xi^2 n}{(U-L)^2}} - 4e^{-\frac{2\zeta^2 n}{(\epsilon^U - \epsilon^L)^2}}$$

By solving ξ to minimize n according to Eq.8 we have:

$$\Pr(|\widehat{I(X, Z)}_{n, \tilde{\theta}} - \widehat{I(X, Z)}_{\infty, \tilde{\theta}}| \leq \epsilon) \geq 1 - \delta. \quad \blacksquare$$

Compared to *MINE*, as per the example shown in §2, for $M = 1$ (i.e. $L = -1$ and $U = 1$), $\delta = 0.05$, $\epsilon = 0.1$, our estimator requires $n = 10,742$ compared to *MINE* requiring $n = 18,756,256$ *i.i.d* validation samples to estimate a lower bound, which makes MI-based dependency analysis feasible for domains where data collection is prohibitively expensive, e.g. fMRI brain scans. In practice, sample complexity can be further optimized by tuning hyperparameters U and L .

Note that the sample complexity of our approach, DEMINE, for estimating Eq.7 does not depend on network size d . The improved sample complexity seemingly comes at a cost of bound tightness

guarantees. In fact, to guarantee bound tightness of Eq.7, $O(d \log d)$ examples would still be theoretically required to learn $\tilde{\theta}$ with guaranteed close values to θ^* , and the total data cost would be on par with *MINE*. In practice, such a learnability bound is known to be overly loose, as over-parameterized neural networks have been shown to generalize well in classification and regression tasks [48]. Fundamentally, what determines bound tightness is the generalization error of $\tilde{\theta}$ – to which the learnability bound is serving as a proxy. Empirically, not only that the bound tightness of DEMINE is as good as *MINE* so the loss of guaranteed tightness did not affect empirical tightness, but the learning-based formulation of DEMINE also allows further bound tightness improvements by learning $\tilde{\theta}$ that generalizes beyond curve fitting using meta-learning.

In the following section, we present a meta-learning formulation, Meta-DEMINE, that learns $\tilde{\theta}$ for generalization given the same model class and training samples.

3.2 Meta-Learning

Given training data $\{(x_i, z_i)_{\text{train}}, i = 1, \dots, m\}$, Meta-DEMINE algorithm first generates MI estimation tasks each consisting of a meta-training split A and a meta-val split B through a novel *task augmentation* process. A parameter initialization θ_{init} is then learned to maximize MI estimation performance on the generated tasks using initialization θ_{init} as shown in Eq.9.

$$\theta_{\text{init}} = \arg \min_{\theta^{(0)} \in \Theta} \mathbb{E}_{(A,B) \in \mathcal{T}} \mathcal{L}((x, z)_{\text{B}}, \theta^{(t)}), \text{ with } \theta^{(t)} \equiv \text{MetaTrain}((x, z)_{\text{A}}, \theta^{(0)}). \quad (9)$$

Here $\theta^{(t)} = \text{MetaTrain}((x, z)_{\text{A}}, \theta^{(0)})$ is the meta-training process of starting from an initialization $\theta^{(0)}$ and applying SGD⁴ over t steps to learn θ where in every meta training iteration we have:

$$\theta^{(t)} \leftarrow \theta^{(t-1)} - \gamma \nabla \mathcal{L}((x, z)_{\text{A}}, \theta^{(t-1)}).$$

Finally, $\tilde{\theta}$ is learned using the entire training set $\{(x_i, z_i)_{\text{train}}, i = 1, \dots, m\}$ with θ_{init} as initialization:

$$\tilde{\theta} = \text{MetaTrain}((x, z)_{\text{train}}, \theta_{\text{init}}).$$

Task Augmentation: Meta-DEMINE adapts MAML [10] for MI lower bound maximization. MAML has been shown to improve generalization performance in N -class K -shot image classification. MI estimation, however, does not come with predefined classes and tasks. A naive approach to produce tasks would be through cross validation – partitioning training data into meta-training and meta-validation splits. However, merely using cross-validation tasks is prone to overfitting – a θ_{init} , which memorizes all training samples would as a result have memorized all meta-validation splits. Instead, Meta-DEMINE generates tasks by augmenting the cross validation tasks through *task augmentation*. Training samples are first split into meta-training and meta-validation splits, and then transformed using the same random invertible transformation to increase task diversity. Meta-DEMINE generates invertible transformation by sequentially composing the following functions:

$$\begin{aligned} \text{Mirror} : & \quad m(x) = (2n - 1)x, & \quad n \sim \text{Bernoulli}(\frac{1}{2}), \\ \text{Permute} : & \quad P(x) = {}^n P_d, & \quad \text{Permute dimensions.} \\ \text{Offset} : & \quad O(x) = x + \epsilon, & \quad \epsilon \sim \mathcal{U}(-0.1, 0.1), \\ \text{Gamma} : & \quad G(x) = \text{sign}(x) |x|^\gamma, & \quad \gamma \sim \mathcal{U}(0.5, 2), \end{aligned}$$

Since the MI between two random variables is invariant to invertible transformations on each variable, MetaTrain is expected to arrive at the same MI lower bound estimation regardless of the transformation applied. At the same time, memorization is greatly suppressed, as the same pair (x, z) can have different $\log \frac{p(x,z)}{p(x)p(z)}$ under different transformations. More sophisticated invertible transformations (affine, piece-wise linear) can also be added. Task augmentation is an orthogonal approach to data augmentation. Using image classification as an example, data augmentation generates variations of the image, translated, or rotated images assuming that they are valid examples of the class. Task augmentation on the other hand, does not make such assumption. Task augmentation requires the initial parameters θ_{init} to be capable of recognizing the same class in a world where all images are translated and/or rotated, with the assumption that the optimal initialization should easily adapt to both the upright world and the translated and/or rotated world.

Optimization: Solving θ_{init} using the meta-learning formulation Eq.9 poses a challenging optimization problem. The commonly used approach is back propagation through time (BPTT) which

⁴In practice, Adam [28] is used for faster optimization. Illustrating SGD for simplicity.

Algorithm 2 Meta-DEMINE

Input Data: $\{(x, z)_{\text{train}}, (x, z)_{\text{val}}\}$
Parameters: batch \mathcal{B} , Meta Learning Iterations N_M , Task Augmentation Iterations N_T , Optimization Iterations N_O , Ratio r , Learning rate η , Meta Learning Rate η_{meta}
Output: MI, $T_{\theta_{\text{init}}}(X, Z)$, $T_{\theta}(X, Z)$

- 1: **for** $i = 1 : N_M$ **do**
- 2: **for** $j = 1 : N_T$ **do**
- 3: $A = r \times \text{train}, B = \text{train} - A$
- 4: Split $(x, z)_{\text{train}}$ into $(x, z)_A$ and $(x, z)_B$
- 5: Transformation R_x for x , $R_x(\cdot) = \text{m}(\text{P}(\text{O}(\text{G}(\cdot))))$
- 6: Transformation R_z for z , $R_z(\cdot) = \text{m}(\text{P}(\text{O}(\text{G}(\cdot))))$
- 7: $\theta_{\text{meta}}^{(0)} \leftarrow \theta_{\text{init}}$
- 8: **for** $k = 1 : N_O$ **do**
- 9: Sample a batch of $(x, z)_B \sim (x, z)_A$
- 10: Compute $\mathcal{L}((R_x(x), R_z(z))_B, \theta_{\text{meta}}^{(k)})$
- 11: Compute $\nabla_{\theta_{\text{meta}}^{(k)}} \mathcal{L}$ – gradient for θ_{meta}
- 12: Update θ_{meta} using Adam [28] with η
- 13: **end for**
- 14: Compute $\mathcal{L}_{\text{meta}}((R_x(x), R_z(z))_B, \theta_{\text{meta}}^{(N_O)})$
- 15: Compute $\nabla_{\theta_0} \mathcal{L}_{\text{meta}}$ – gradient to θ_{init} using BPTT
- 16: **end for**
- 17: Update θ_{init} using Adam [28] with η_{meta}
- 18: **end for**
- 19: $\theta^{(0)} \leftarrow \theta_{\text{init}}$
- 20: **for** $i = 1 : N_O$ **do**
- 21: Sample a batch of $(x, z)_B \sim (x, z)_{\text{train}}$
- 22: Compute $\mathcal{L}((x, z)_B, \theta^{(i)})$
- 23: Compute gradient $\nabla_{\theta} \mathcal{L}$
- 24: Update θ using Adam with η
- 25: **end for**
- 26: Compute MI = $\mathcal{L}((x, z)_{\text{val}}, \theta^{(N_O)})$
- 27: **return** MI, $\theta_{\text{init}}, \theta^{(N_O)}$

computes second order gradients and directly back propagate gradient from $\text{MetaTrain}((x, z)_A, \theta^{(0)})$ to θ_{init} . BPTT is very effective for a small number of optimization steps, but is vulnerable to exploding gradients and is memory intensive. In addition to BPTT, we find that stochastic finite difference algorithms such as Evolution Strategies (ES) [37] and Parameter-Exploring Policy Gradients (PEPG) [39] can sometimes improve optimization. In practice, we use BPTT or PEPG to optimize Eq.9 depending on the problem. Meta-DEMINE algorithm is specified in Algorithm 2.

4 Evaluation on Synthetic Datasets

Dataset. We evaluate our approaches DEMINE and Meta-DEMINE against baselines and state-of-the-art approaches on 3 synthetic datasets: 1D Gaussian, 20D Gaussian and sine wave. For 1D and 20D Gaussian datasets, following [5], we define two k -dimensional multivariate Gaussian random variables X and Z which have component-wise correlation $\text{corr}(X_i, Z_j) = \delta_{ij}\rho$, where $\rho \in (-1, 1)$ and δ_{ij} is Kronecker’s delta. Mutual information $I(X; Z)$ has a closed form solution $I(X; Z) = -k \ln(1 - \rho^2)$. For sine wave dataset, we define two random variables X and Z , where $X \sim \mathcal{U}(-1, 1)$, $Z = \sin(aX + \frac{\pi}{2}) + 0.05\epsilon$, and $\epsilon \sim \mathcal{N}(0, 1)$. Estimating mutual information accurately given few pairs of (X, Z) requires the ability to extrapolate the sine wave given few examples. Ground truth MI for sine wave dataset is approximated by running the the KSG Estimator [29] on 1,000,000 samples.

Implementation. We compare our estimators, DEMINE and Meta-DEMINE, against the KSG estimator [29] MI-KSG and MINE-f. For both DEMINE and Meta-DEMINE, we study variance reduction mode, referred to as *-vr*, where hyperparameters are selected by optimizing 95% confident estimation mean $(\mu - 2\sigma_{\mu})$ and statistical significance mode, referred to as *-sig*, where hyperparameters

are selected by optimizing 95% confident MI lower bound ($\mu - \epsilon$). Samples (x, z) are split into 50%-50% as $(x, z)_{\text{train}}$ and $(x, z)_{\text{val}}$.

We use a separable network architecture $T_\theta(x, z) = M(\tanh(w \cos \langle f(x), g(z) \rangle) + b) - t$. f and g are MLP encoders that embed signals x and z into vector embeddings. Hyperparameters $t \in [-1, 1]$ and M control upper and lower bounds $T_\theta(x, z) \in [-M(1+t), M(1-t)]$. Parameters w and b are learnable parameters. MLP design and optimization hyperparameters are selected using Bayesian hyperparameter optimization [6] with 3-fold cross-validation on $(x, z)_{\text{train}}$ over 1,000 iterations.

Hyperparameter search on DEMINE-vr and DEMINE-sig was conducted using the hyperopt package⁵. Seven hyper parameters were involved in hyperparameter search: 1) number of encoder layers [1, 5], 2) encoder hidden size [8, 256], 3) learning rate η [10^{-4} , 3×10^{-1}] in log scale, 4) number of optimization iterations N_O [5, 200] (sine wave [5, 5000]) in log scale, 5) batch size \mathcal{B} [256, 1024], 6) M , [10^{-3} , 5] in log scale, 7) t , [-1, 1]. Mean μ and sample standard deviation σ of MI estimate computed over 3 fold cross validation on $(x, z)_{\text{train}}$. DEMINE-vr maximizes two sigma low $\mu - 2\sigma_\mu$ where $\sigma_\mu = \frac{1}{\sqrt{3}}\sigma$. DEMINE-sig maximizes statistical significance $\mu - \epsilon$ where ϵ is two-sided 95% confidence interval of MI. Meta-DEMINE-vr and Meta-DEMINE-sig subsequently reuse these hyperparameters as DEMINE-vr and DEMINE-sig.

Meta-learning hyperparameters are chosen as outer loop $N_M = 3,000$ iterations, task augmentation $N_T = 1$ iterations, $r = 0.8$, $\eta_{\text{meta}} = \frac{\eta}{3}$, with task augmentation mode $m(P(O(\cdot)))$. N_O capped at 30 iterations for 1D and 20D Gaussian datasets due to memory limit. The sine wave datasets require large N_O , we used PEPG [39] rather than BPTT.

For MI-KSG, we use off-the-shelf implementation [15] with default number of nearest neighbors $k=3$. MI-KSG does not provide any confidence interval. For MINE-f, we use the same network architecture same as DEMINE-vr. we implement both the original formulation which optimizes T_θ on (x, z) till convergence (10k iters), as well as our own implementation MINE-f-ES with early stopping, where optimization is stopped after the same number of iterations as DEMINE-vr to control overfitting.

Results. Figure 1(a) shows MI estimation performance on 20D Gaussian datasets with varying $\rho \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$ using $N = 300$ samples. Results are averaged over 5 runs to compare estimator bias, variance and confidence. Note that Meta-DEMINE-sig detects the highest $p < 0.05$ confidence MI, outperforming DEMINE-sig which is a close second. Both detect $p < 0.05$ statistically significant dependency starting $\rho = 0.3$, whereas estimations of all other approaches are low confidence. It shows that in contrary to common belief, estimating the variational lower bounds with high confidence can be challenging under limited data. MINE-f estimates $\text{MI} > 3.0$ and MINE-f-ES estimates positive MI when $\rho = 0$, both due to overfitting, despite MINE-f-ES having the lowest empirical bias. DEMINE variants have relatively high empirical bias but low variance due to tight upper and lower bound control, which provides a different angle to understand bias-variance trade off in MI estimation [35].

Figure 1(b,c,d) shows MI estimation performance on 1D, 20D Gaussian and sine wave datasets with fixed $\rho = 0.8, 0.3$ and $a = 8\pi$ respectively, with varying $N \in \{30, 100, 300, 1000, 3000\}$ number of samples. More samples asymptotically improves empirical bias across all estimators. As opposed to 1D Gaussian datasets which are well solved by $N = 300$ samples, higher-dimensional 20D Gaussian and higher-complexity sine wave datasets are much more challenging and are not solved using $N = 3000$ samples with a signal-agnostic MLP architecture. DEMINE-sig and Meta-DEMINE-sig detect $p < 0.05$ statistically significant dependency on not only 1D and 20D Gaussian datasets where x and z have non-zero correlation, but also on the sine wave datasets where correlation between x and z is 0. This means that DEMINE-sig and Meta-DEMINE-sig can be used for nonlinear dependency testing to complement linear correlation testing.

We study the effect of cross-validation meta-learning and task augmentation on 20D Gaussian with $\rho = 0.3$ and $N = 300$. Figure 2 plots performance of Meta-DEMINE-vr over $N_M = 3000$ meta iterations under combinations of task augmentations modes and number of adaptation iterations $N_O \in \{0, 20\}$. Overall, task augmentation modes which involve axis flipping $m(\cdot)$ and permutation $P(\cdot)$ are the most successful. With $N_O = 20$ steps of adaptation, task augmentation modes $P(\cdot)$, $m(P(\cdot))$ and $m(P(O(\cdot)))$ prevent overfitting and improves performance. The performance improvements of task augmentation is not simply from change in batch size, learning rate or number of optimization

⁵Hyperopt package: <https://github.com/hyperopt/hyperopt>.

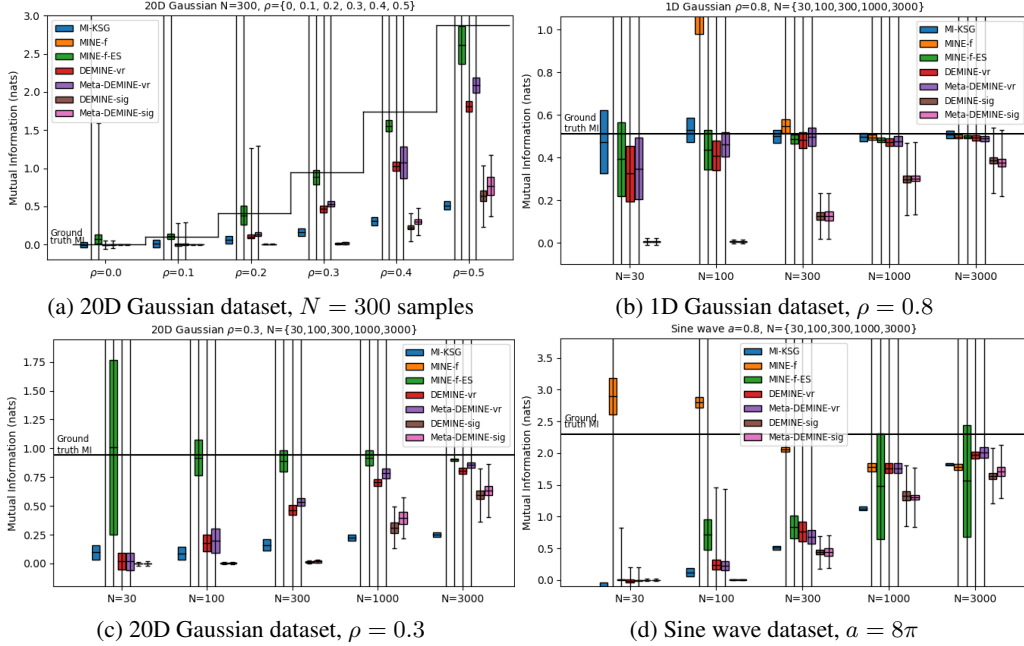


Figure 1: Comparing MI Estimation performance of DEMINE and Meta-DEMINE with the KSG estimator [29] and MINE-f [5] on different datasets using varying number of samples. The bars show estimator mean and standard deviation averaged over 5 runs with different seeds. The errorbars show 95% confidence interval (not available for MI-KSG). The statistical significance focused variants DEMINE-sig and Meta-DEMINE-sig achieves the highest 95% confident MI estimation. Meta-DEMINE improves over DEMINE most of the time. Best viewed in color.

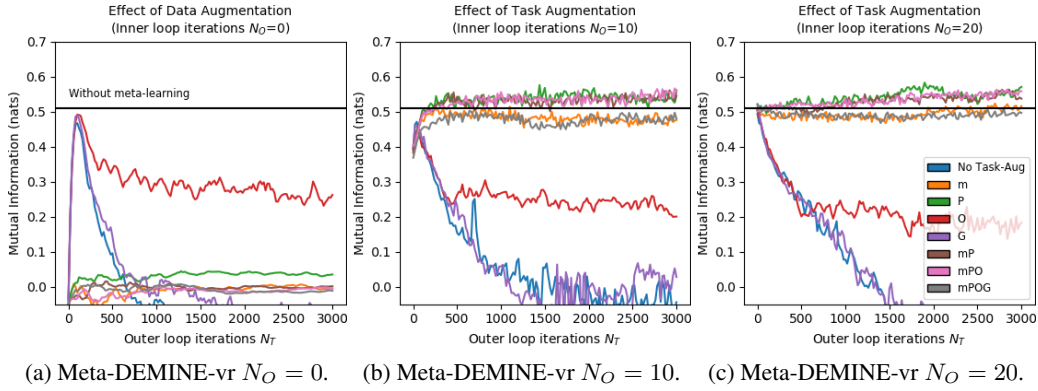


Figure 2: To study the effect of task augmentation and number of adaptation steps, we run Meta-DEMINE-vr with different task augmentation modes and vary number of adaptation iterations $N_O \in \{0, 10, 20\}$ on Gaussian 20D, $\rho = 0.3$ dataset. Combinations of permutation and mirroring operations are effective in reducing overfitting and improving performance. Best viewed in color.

Table 1: Number of HCP-MMP1 regions with significant pairwise correlation (r) and MI (DEMINE, Meta-DEMINE) during listening.

No. shared	r	DEMINE	Meta-DEMINE
r	37	24	23
DEMINE	24	28	26
Meta-DEMINE	23	26	29

Table 2: Segment classification accuracy for NeuralMI versus Pearson’s correlation in 1-vs-1 and 1-vs-rest*.

Classification Accuracy (%)	ISC Mask					dDMN Mask				
	P	F	Br	Bk	MI	P	F	Br	Bk	MI
Chance	3.7	1.8	2.6	1.9	N/A	3.7	1.8	2.6	1.9	N/A
Pearson’s r 1vR	35.0	20.4	25.8	31.5	N/A	14.8	6.4	11.8	9.9	N/A
DEMINE 1vR	42.8	28.0	32.8	35.9	0.637	16.5	7.9	11.6	12.0	0.035
Meta-DEMINE 1vR	47.2	32.5	39.9	41.0	0.752	13.7	7.9	8.2	8.9	0.031

Abbreviations: P: Pieman; F: Forgotten; Br: Bronx; Bk: Black, MI: Mutual Information.
*Note that all the results are averaging over other subjects.

iterations, because meta-learning without task augmentation for both $N_O = 0$ and 20 could not outperform baseline. Meta-learning without task augmentation and with task augmentation but using only $O(\cdot)$ or $G(\cdot)$ result in overfitting. Task augmentation with $m(\cdot)$ or $m(P(O(G(\cdot))))$ prevent overfitting, but do not provide performance benefits, possibly because their complexity is insufficient or excessive for 20 adaptation steps. Further more, task augmentation with no adaptation ($N_O = 0$) falls back to data augmentation, where samples from transformed distributions are directly used to learn $T_\theta(x, z)$. Data augmentation with $O(\cdot)$ outperforms no augmentation, but is unable to outperform baseline and suffer from overfitting. It shows that task augmentation provides improvements orthogonal to data augmentation.

5 Application: fMRI Inter-subject correlation (ISC) analysis

Humans use language to effectively transmit brain representations among conspecifics. For example, after witnessing an event in the world, a speaker may use verbal communication to evoke neural representations reflecting that event in a listener’s brain [23]. The efficacy of this transmission, in terms of listener comprehension, is predicted by speaker–listener neural synchrony and synchrony among listeners [43]. To date, most work has measured brain-to-brain synchrony by locating statistically significant inter-subject correlation (ISC); quantified as the Pearson product-moment correlation coefficient between response time series for corresponding voxels or regions of interest (ROIs) across individuals [24, 38, 40]. Using DEMINE and Meta-DEMINE for statistical dependency testing, we can extend ISC analysis to capture nonlinear and higher-order interactions in continuous fMRI responses. Specifically, given synchronized fMRI response frames in two brain regions X and Z across K subjects $X_i, Z_i, i = 1, \dots, K$ as random variables. We model the conditional mutual information $I(X_i; Z_j | i \neq j)$ as the MI form of pair-wise ISC analysis. By definition, $I(X_i; Z_j | i \neq j)$ first computes MI between activations X_i and Z_j from subjects i and j respectively, and then average across pairs of subjects $i \neq j$. It can be lower bounded using Eq. 7 by learning a $T_\theta(x, z)$ shared across all subject pairs.

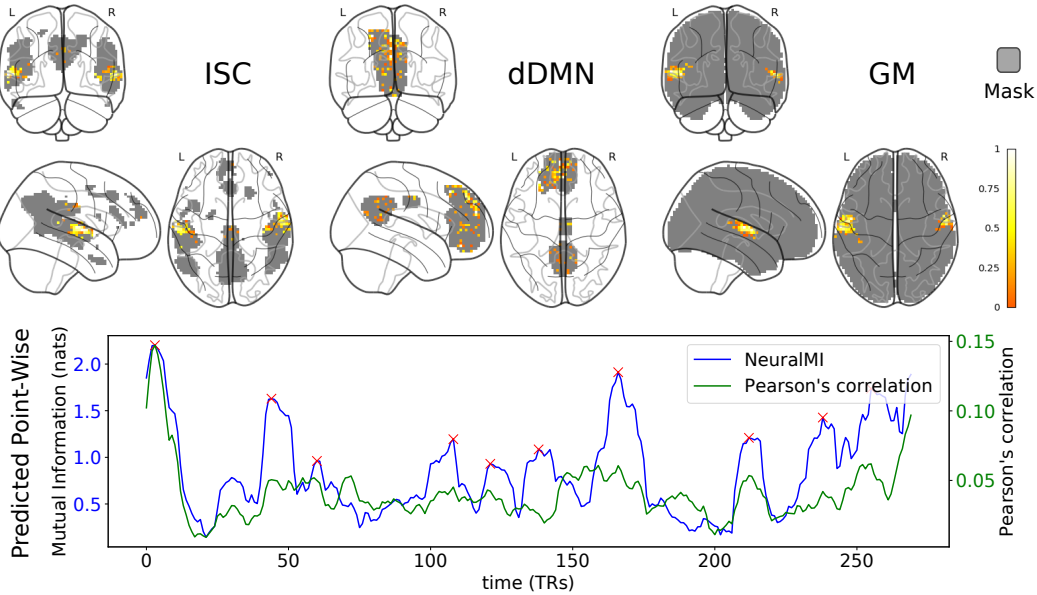
Dataset. We study MI-based and correlation-based ISC on a fMRI story comprehension dataset [41] with 40 participants listening to four spoken stories. Average story duration is 11 minutes. An fMRI frame with full brain coverage is captured at repetition time 1 TR = 1.5 seconds with 2.5mm isotropic spatial resolution. We restricted our analysis to subsets of voxels defined using independent data from previous studies: functionally-defined masks of high ISC voxels (ISC; 3,800 voxels) and dorsal Default-Mode Network voxels (dDMN; 3,940 voxels) from [41] as well as 180 HCP-MMP1 multimodal cortex parcels from [17]. All masks were defined in MNI space.

Implementation. We compare MI-based ISC using DEMINE and Meta-DEMINE with correlation-based ISC using Pearson’s correlation. DEMINE and Meta-DEMINE setup follows Section §4. The fMRI data were partitioned by subject into a train set of 20 subjects and a validation set of 20 different subjects. Residual 1D CNN is used instead of MLP as the encoder for studying temporal dependency. For Pearson’s correlation, high-dimensional signals are reshaped to 1D for correlation analysis.

Quantitative Results. We first study that for the fine grained HCM-MMP1 brain regions, which of them have $p < 0.05$ statistically significant activities by MI and Pearson’s correlation. Table 1 shows the result. Overall, more regions have statistically significant correlation than dependency. This is expected because correlation requires less data to detect. But Meta-DEMINE is able to find 6 brain regions that potentially have statistically significant dependency but lacks significant correlation. This shows that MI analysis can be used to complement correlation-based ISC analysis.

By considering temporal ISC over time, fMRI signals can be modeled with improved accuracy. In Table 2 we apply DEMINE and Meta-DEMINE with $L = 10$ TRs (15s) sliding windows as random variables to study amount of information that can be extracted from ISC and dDMN masks. We use between-subject time-segment classification (BSC) for evaluation [25, 22]. Each fMRI scan is divided into K non-overlapping $L = 10$ TRs time segments. The BSC task is one versus rest retrieval: retrieve the corresponding time segment z of an individual given a group of time segments x excluding that individual, measured by top-1 accuracy. For retrieval score, $T_\theta(X, Z)$ is used for DEMINE and Meta-DEMINE and $\rho(X, Z)$ is used for Pearson’s correlation as a simple baseline. With CNN as encoder, DEMINE and Meta-DEMINE model the signal better and achieve higher accuracy. Also, Meta-DEMINE is able to extract 0.75 nats of MI from the ISC mask over 10TRs or 15s, which could potentially be improved by more samples and high frequency fMRI scans.

Figure 3: Top: Top contributing voxels in the learned $T_\theta(X, Z)$ by gradient magnitude $\mathbb{E}_X(\frac{\partial T}{\partial X_i})^2$. Auditory region is highlighted for ISC and GM masks (best in color). Bottom: Evaluation on the "Pie Man" dataset using the ISC mask showing our approach $T_\theta(X, Z)$ versus Pearson correlation over time in the one versus rest case averaged over 20 test subjects.



Qualitative Results. Fig. 3 (top) visualizes voxels that are important to $T_\theta(x, z)$ of the DEMINE model using their gradient magnitude variance for the ISC and dDMN masks, as well as an anatomically-defined Gray Matter (GM) mask. The DEMINE model focuses on auditory regions functionally important for perceiving the story stimulus.

Fig. 3 (bottom) plots the $T(x, z)$ and inter-subject Pearson correlations over time for "Pie Man" using the ISC mask and a sliding window size $L = 10$, using the one vs rest scores averaged over all subjects. DEMINE yields more distinctive peaks.

We identify the peaks in DEMINE for "Pie Man" (with Pearson correlations) over time, then locate the story transcriptions in the $L = 10TRs$ (15 seconds) window corresponding to the peak:

- 4: "... toiled for The Ram, uh, Fordham University's student newspaper. And one day, I'm walking toward the campus center and out comes the elusive Dean McGowen, architect of a policy to replace traditionally ..."
- 45: "The Dean is covered with cream. So I give him a moment, then I say, 'Dean McGowen, would you care to comment on this latest attack?' And he says, 'Yes, I would care to comment. ...'"
- 109: "... which makes no sense. Fordham was a Catholic school and we all thought Latin was classy so, that's what I used. And when I finished my story, I, I raced back to Dwyer and I showed it to him and he read it and he said ..."
- 122: "Few days later, I get a letter. I opened it up and it says, 'Dear Jim, good story. Nice details. If you want to see me again in action, be on the steps of Duane Library ...'"
- 139: "... out comes student body president, Sheila Biel. And now, Sheila Biel was different from the rest of us flannel-shirt wearing, part-time-job working, Fordham students. Sheila was..."
- 167: "Pie Man emerged from behind a late night library drop box, made his delivery, and fled away, crying, 'Ego sum non una bestia.' And that's what I reported in my story..."
- 213: "... that there was a question about whether she even knew if I existed. So I saw her there and made a mental note to do nothing about it, and then I went to the bar and ordered a drink, and I felt a, a tap on my shoulder. I turned around, and it was her..."

- 239: "And wasn't I really Pie Man? Hadn't I brought him into existence? Didn't she only know about him because of me? But actually ..."
- 256: "I said, 'Yes, Angela, I am Pie Man.' And she looked at me and she said, 'Oh, good. I was hoping you'd say that ...'"

We hypothesize that the scripts associated with the peaks may capture points when listeners pay more attention, resulting in the Signal-to-Noise Ratio (SNR) of fMRI scans being enhanced.

6 Conclusion

We illustrated that a predictive view of the MI lower bounds coupled with meta-learning results in data-efficient variational MI estimators, DEMINE and Meta-DEMINE, that are capable of performing statistical test of dependency. We also showed that our proposed task augmentation reduces overfitting and improves generalization in meta-learning. We successfully applied MI estimation to real world, data scarce, fMRI datasets. Our results suggest a greater avenue of using neural networks and meta-learning to improve MI analysis and applying neural network-based information theory tools to enhance the analysis of information processing in the brain. Model-agnostic, high-confidence, MI lower bound estimation approaches – including *MINE*, DEMINE and Meta-DEMINE– are limited to estimating small MI lower bounds up to $O(\log N)$ as pointed out in [31], where N is the number of samples. In real fMRI datasets, however, strong dependency is rare and existing MI estimation tools are limited more by their ability to accurately characterize the dependency. Nevertheless, when quantitatively measuring strong dependency, cross-entropy [31] or model-based quantities, alternatives to MI, such as correlation or CCA, may be measured with high confidence.

Acknowledgments

This work is funded by DARPA FA8750-18-C-0213. The views, opinions, and/or conclusions contained in this paper are those of the author and should not be interpreted as representing the official views or policies, either expressed or implied of the DARPA or the DoD.

References

- [1] D. B. F. Agakov. The IM algorithm: a variational approach to information maximization. *Advances in Neural Information Processing Systems*, 16:201, 2004.
- [2] I. Ahmad and P.-E. Lin. A nonparametric estimation of the entropy for absolutely continuous distributions (corresp.). *IEEE Transactions on Information Theory*, 22(3):372–375, 1976.
- [3] B. B. Avants, C. L. Epstein, M. Grossman, and J. C. Gee. Symmetric diffeomorphic image registration with cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain. *Medical Image Analysis*, 12(1):26–41, 2008.
- [4] Y. Behzadi, K. Restom, J. Liau, and T. T. Liu. A component based noise correction method (CompCor) for BOLD and perfusion based fMRI. *NeuroImage*, 37(1):90–101, 2007.
- [5] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, D. Hjelm, and A. Courville. Mutual information neural estimation. In *International Conference on Machine Learning*, pages 530–539, 2018.
- [6] J. Bergstra, D. Yamins, and D. D. Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. 2013.
- [7] R. W. Cox. AFNI: software for analysis and visualization of functional magnetic resonance neuroimages. *Computers and Biomedical research*, 29(3):162–173, 1996.
- [8] C. Daniel. I knew you were black. <https://themoth.org/stories/i-knew-you-were-black>, 2018. Accessed: 2018-10-12.
- [9] O. Esteban, C. Markiewicz, R. W. Blair, C. Moodie, A. I. Isik, A. Erramuzpe Aliaga, J. Kent, M. Goncalves, E. DuPre, M. Snyder, H. Oya, S. Ghosh, J. Wright, J. Durnez, R. Poldrack, and K. J. Gorgolewski. FMRIPrep: a robust preprocessing pipeline for functional MRI. *bioRxiv*, 2018.
- [10] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 1126–1135, 2017.
- [11] C. Finn, K. Xu, and S. Levine. Probabilistic model-agnostic meta-learning. In *Advances in Neural Information Processing Systems*, pages 9537–9548, 2018.
- [12] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine. One-shot visual imitation learning via meta-learning. In *Conference on Robot Learning*, pages 357–368, 2017.
- [13] V. S. Fonov, A. C. Evans, R. C. McKinstry, C. Almlil, and D. Collins. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. *NeuroImage*, (47):S102, 2009.
- [14] N. Gaiman. The man who forgot ray bradbury. <https://soundcloud.com/neilgaiman/the-man-who-forgot-ray-bradbury>, 2018. Accessed: 2018-10-12.
- [15] W. Gao, S. Kannan, S. Oh, and P. Viswanath. Estimating mutual information for discrete-continuous mixtures. In *Advances in Neural Information Processing Systems*, pages 5986–5997, 2017.
- [16] W. Gao, S. Oh, and P. Viswanath. Demystifying fixed k -nearest neighbor information estimators. *IEEE Transactions on Information Theory*, 64(8):5629–5661, 2018.
- [17] M. F. Glasser, T. S. Coalson, E. C. Robinson, C. D. Hacker, J. Harwell, E. Yacoub, K. Ugurbil, J. Andersson, C. F. Beckmann, M. Jenkinson, et al. A multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171, 2016.
- [18] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS’10). Society for Artificial Intelligence and Statistics*, 2010.
- [19] K. Gorgolewski, C. Burns, C. Madison, D. Clark, Y. Halchenko, M. Waskom, and S. Ghosh. Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in python. *Frontiers in Neuroinformatics*, 5:13, 2011.
- [20] K. J. Gorgolewski, T. Auer, V. D. Calhoun, R. C. Craddock, S. Das, E. P. Duff, G. Flandin, S. S. Ghosh, T. Glatard, Y. O. Halchenko, et al. The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Scientific Data*, 3:160044, 2016.
- [21] D. N. Greve and B. Fischl. Accurate and robust brain image alignment using boundary-based registration. *NeuroImage*, 48(1):63–72, 2009.
- [22] J. S. Guntupalli, M. Hanke, Y. O. Halchenko, A. C. Connolly, P. J. Ramadge, and J. V. Haxby. A model of representational spaces in human cortex. *Cerebral Cortex*, 26(6):2919–2934, 2016.

- [23] U. Hasson, A. A. Ghazanfar, B. Galantucci, S. Garrod, and C. Keysers. Brain-to-brain coupling: a mechanism for creating and sharing a social world. *Trends in cognitive sciences*, 16(2):114–121, 2012.
- [24] U. Hasson, Y. Nir, I. Levy, G. Fuhrmann, and R. Malach. Intersubject synchronization of cortical activity during natural vision. *Science*, 303(5664):1634–1640, 2004.
- [25] J. V. Haxby, J. S. Guntupalli, A. C. Connolly, Y. O. Halchenko, B. R. Conroy, M. I. Gobbini, M. Hanke, and P. J. Ramadge. A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron*, 72(2):404–416, 2011.
- [26] M. Jenkinson, P. Bannister, M. Brady, and S. Smith. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *NeuroImage*, 17(2):825–841, 2002.
- [27] T. Kim, J. Yoon, O. Dia, S. Kim, Y. Bengio, and S. Ahn. Bayesian model-agnostic meta-learning. *arXiv preprint arXiv:1806.03836*, 2018.
- [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [29] A. Kraskov, H. Stogbauer, and P. Grassberger. Estimating mutual information. *Physical review E*, 2004.
- [30] D. Maclaurin, D. Duvenaud, and R. Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pages 2113–2122, 2015.
- [31] D. McAllester and K. Statos. Formal limitations on the measurement of mutual information. *arXiv preprint arXiv:1811.04251*, 2018.
- [32] J. O’Grady. Pie Man. <https://themoth.org/stories/pie-man>, 2018. Accessed: 2018-10-12.
- [33] J. O’Grady. Running from the Bronx. <https://soundcloud.com/the-story-collider/jim-ogrady-running-from-the>, 2018. Accessed: 2018-10-12.
- [34] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean. Efficient neural architecture search via parameter sharing. In *International Conference on Machine Learning*, pages 4092–4101, 2018.
- [35] B. Poole, S. Ozair, A. van den Oord, A. A. Alemi, and G. Tucker. On variational lower bounds of mutual information. In *Bayesian Deep Learning Workshop, NeurIPSW*, 2018.
- [36] J. D. Power, A. Mitra, T. O. Laumann, A. Z. Snyder, B. L. Schlaggar, and S. E. Petersen. Methods to detect, characterize, and remove motion artifact in resting state fMRI. *NeuroImage*, 84:320–341, 2014.
- [37] T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- [38] M. B. Schippers, A. Roebroeck, R. Renken, L. Nanetti, and C. Keysers. Mapping the information flow from one brain to another during gestural communication. *Proceedings of the National Academy of Sciences*, page 201001791, 2010.
- [39] F. Sehnke, C. Osendorfer, T. Rückstieß, A. Graves, J. Peters, and J. Schmidhuber. Parameter-exploring policy gradients. *Neural Networks*, 23(4):551–559, 2010.
- [40] L. J. Silbert, C. J. Honey, E. Simony, D. Poeppel, and U. Hasson. Coupled neural systems underlie the production and comprehension of naturalistic narrative speech. *Proceedings of the National Academy of Sciences*, 111(43):E4687–E4696, 2014.
- [41] E. Simony, C. J. Honey, J. Chen, O. Lositsky, Y. Yeshurun, A. Wiesel, and U. Hasson. Dynamic reconfiguration of the default mode network during narrative comprehension. *Nature Communications*, 7:12141, 2016.
- [42] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 4077–4087, 2017.
- [43] G. J. Stephens, L. J. Silbert, and U. Hasson. Speaker–listener neural coupling underlies successful communication. *Proceedings of the National Academy of Sciences*, 107(32):14425–14430, 2010.
- [44] J. M. Treiber, N. S. White, T. C. Steed, H. Bartsch, D. Holland, N. Farid, C. R. McDonald, B. S. Carter, A. M. Dale, and C. C. Chen. Characterization and correction of geometric distortions in 814 diffusion weighted images. *PLOS ONE*, 11(3):e0152472, 2016.
- [45] N. J. Tustison, B. B. Avants, P. A. Cook, Y. Zheng, A. Egan, P. A. Yushkevich, and J. C. Gee. N4itk: improved n3 bias correction. *IEEE Transactions on Medical Imaging*, 29(6):1310–1320, June 2010.

- [46] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016.
- [47] S. Wang, D. J. Peterson, J. C. Gatenby, W. Li, T. J. Grabowski, and T. M. Madhyastha. Evaluation of field map and nonlinear registration methods for correction of susceptibility artifacts in diffusion mri. *Frontiers in Neuroinformatics*, 11:17, 2017.
- [48] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [49] Y. Zhang, M. Brady, and S. Smith. Segmentation of brain MR images through a hidden markov random field model and the expectation-maximization algorithm. *IEEE Transactions on Medical Imaging*, 20(1):45–57, 2001.

A Additional Details about the fMRI Dataset

The dataset we use [41], contains 40 participants (mean age = 23.3 years, SD = 8.9, range: 18–53; 27 female) recruited to listen to four spoken stories⁶⁷. The stories were renditions of “Pie Man” and “Running from the Bronx” by Jim O’Grady [32, 33], “The Man Who Forgot Ray Bradbury” by Neil Gaiman [14], and “I Knew You Were Black” by Carol Daniel [8]; story durations were 7, 9, 14, and 13 minutes, respectively. After scanning, participants completed a questionnaire comprising 25–30 questions per story intended to measure narrative comprehension. The questionnaires included multiple choice, True/False, and fill-in-the-blank questions, as well as four additional subjective ratings per story. Functional and structural images were acquired using a 3T Siemens Prisma with a 64-channel head coil (see Section §A.1 for additional details). Briefly, functional images were acquired in an interleaved fashion using gradient-echo echo-planar imaging with a multiband acceleration factor of 3 (TR/TE = 1500/31 ms, resolution = 2.5 mm isotropic voxels, full brain coverage).

All fMRI data were formatted according to the Brain Imaging Data Structure (BIDS) standard [20] and preprocessed using fMRIPrep [9] (see Section §A.2 for additional details). Functional data were corrected for slice timing, head motion, and susceptibility distortion, and normalized to MNI space using nonlinear registration. Nuisance variables comprising head motion parameters, framewise displacement, linear and quadratic trends, sine/cosine bases for high-pass filtering (0.007 Hz), and six principal component time series from cerebrospinal fluid (CSF) and white matter were regressed out of the signal using AFNI [7].

The fMRI data comprise $\mathcal{X} \in \mathbb{R}^{V_i \times T}$ for each subject, where V_i represents the flattened and masked voxel space and T represents the number of samples (TRs) during auditory stimulus presentation.

A.1 Additional Details on Dataset Collection

Functional and structural images were acquired using a 3T Siemens Magnetom Prisma with a 64-channel head coil. Functional, blood-oxygenation-level-dependent (BOLD) images were acquired in an interleaved fashion using gradient-echo echo-planar imaging with pre-scan normalization, fat suppression, a multiband acceleration factor of 3, and no in-plane acceleration: TR/TE = 1500/31 ms, flip angle = 67°, bandwidth = 2480 Hz/Px, resolution = 2.5 mm³ isotropic voxels, matrix size = 96 x 96, FoV = 240 x 240 mm, 48 axial slices with roughly full brain coverage and no gap, anterior–posterior phase encoding. At the beginning of each scanning session, a T1-weighted structural scan was acquired using a high-resolution single-shot MPRAGE sequence with an in-plane acceleration factor of 2 using GRAPPA: TR/TE/TI = 2530/3.3/1100 ms, flip angle = 7°, resolution = 1.0 x 1.0 x 1.0 mm voxels, matrix size = 256 x 256, FoV = 256 x 256 x 176 mm, 176 sagittal slices, ascending acquisition, anterior–posterior phase encoding, no fat suppression, 5 min 53 s total acquisition time. At the end of each scanning session a T2-weighted structural scan was acquired using the same acquisition parameters and geometry as the T1-weighted structural image: TR/TE = 3200/428 ms, 4 min 40 s total acquisition time. A field map was acquired at the beginning of each scanning session, but was not used in subsequent analyses.

A.2 Additional Details on Dataset Preprocessing

Preprocessing was performed using fMRIPrep [9], a Nipype [19] based tool. T1-weighted images were corrected for intensity non-uniformity using N4 bias field correction [45] and skull-stripped using ANTs [3]. Nonlinear spatial normalization to the ICBM 152 Nonlinear Asymmetrical template version 2009c [13] was performed using ANTs. Brain tissue segmentation cerebrospinal fluid (CSF), white matter, and gray matter was performed using FSL’s FAST [49]. Functional images were slice timing corrected using AFNI’s 3dTshift [7] and corrected for head motion using FSL’s MCFLIRT [26]. “Fieldmap-less” distortion correction was performed by co-registering each subject’s functional image to that subject’s intensity-inverted T1-weighted image [47] constrained with an average field map template [44]. This was followed by co-registration to the corresponding T1-weighted image using FreeSurfer’s boundary-based registration [21] with 9 degrees of freedom. Motion correcting transformations, field distortion correcting warp, BOLD-to-T1 transformation and T1-to-template (MNI) warp were concatenated and applied in a single step with Lanczos interpolation using ANTs.

⁶Two of the stories were told by a professional storyteller undergoing an fMRI scan; however, fMRI data for the speaker were not analyzed for the present work due to the head motion induced by speech production.

⁷The study was conducted in compliance with the Institutional Review Board of the University

Physiological noise regressors were extracted applying aCompCor [4]. Six principal component time series were calculated within the intersection of the subcortical mask and the union of CSF and WM masks calculated in T1w space, after their projection to the native space of each functional run. Framewise displacement [36] was calculated for each functional run. Functional images were downsampled to 3 mm resolution. Nuisance variables comprising six head motion parameters (and their derivatives), framewise displacement, linear and quadratic trends, sine/cosine bases for high-pass filtering (0.007 Hz cutoff), and six principal component time series from an anatomically-defined mask of cerebrospinal fluid (CSF) and white matter were regressed out of the signal using AFNI's 3dTproject [7]. Functional response time series were z-scored for each voxel.