

ABSTRACT

- We propose a new way of constructing invertible neural networks by combining simple building blocks with a novel set of composition rules, which leads to a rich set of invertible architectures, including those similar to ResNets.
- Inversion is achieved with a locally convergent iterative procedure that is parallelizable and very fast in practice. Additionally, the determinant of the Jacobian can be computed analytically and efficiently, enabling their generative use as flow models.
- We show that our invertible neural networks are competitive with ResNets on MNIST and CIFAR-10 classification. When trained as generative models, our invertible networks achieve likelihoods comparable to state-of-the-art results on MNIST, CIFAR-10 and ImageNet 32×32

Code:



BUILDING INVERTIBLE MODULES

The basic module:

$f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$, with $\mathbf{W} \in \mathbb{R}^{D \times D}$, and $\mathbf{b} \in \mathbb{R}^D$. \mathbf{W} is a triangular matrix.

The calculus of building invertible modules:

Proposition 1: Define \mathcal{F} as the set of all continuously differentiable functions whose Jacobian is lower triangular. Then \mathcal{F} contains the basic module, and is closed under the following composition rules.

- Rule of addition:** $f_1 \in \mathcal{F} \wedge f_2 \in \mathcal{F} \Rightarrow \lambda f_1 + \mu f_2 \in \mathcal{F}$, where $\lambda, \mu \in \mathbb{R}$.
- Rule of composition:** $f_1 \in \mathcal{F} \wedge f_2 \in \mathcal{F} \Rightarrow f_2 \circ f_1 \in \mathcal{F}$. A special case is $f \in \mathcal{F} \Rightarrow h \circ f \in \mathcal{F}$, where $h(\cdot)$ is a continuously differentiable non-linear activation function that is applied element-wise.

Theorem 1: If $J_f \in \mathcal{F}$ is non-singular for all \mathbf{x} in the domain, then f is invertible.

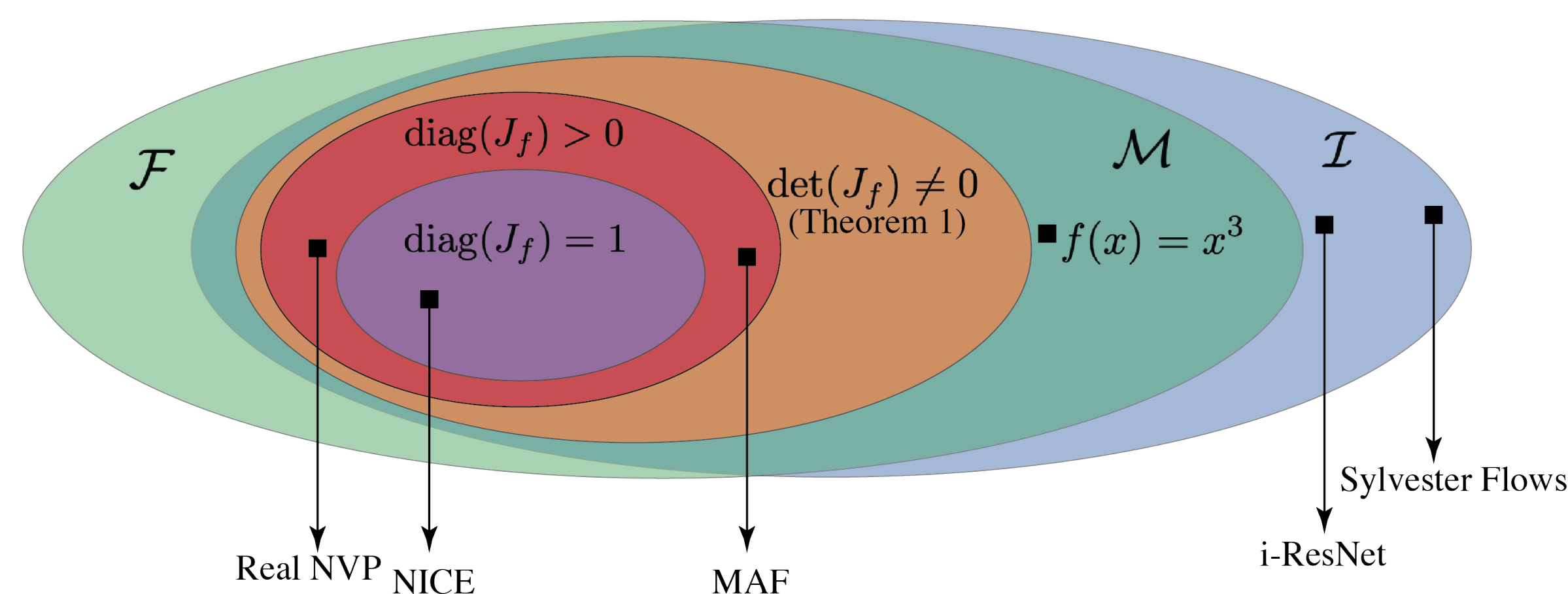


Figure 1: Venn Diagram relationships between invertible functions (\mathcal{F}).

Efficient and parallelizable inversion of the modules:

T : number of iteration steps

α : step size

Algorithm 1: Fixed-point integration for inversion

Require: T, α

- 1: Initialize \mathbf{x}_0
- 2: **for** $t \leftarrow 1$ to T **do**
- 3: Compute $f(\mathbf{x}_{t-1})$
- 4: Compute $\text{diag}(J_f(\mathbf{x}_{t-1}))$
- 5: $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \alpha \text{diag}(J_f(\mathbf{x}_{t-1}))^{-1}(f(\mathbf{x}_{t-1}) - \mathbf{z})$
- 6: **end for**

return \mathbf{x}_T

Theorem 2: The fixed-point iteration method is locally convergent whenever $0 < \alpha < 2$.

MASKED INVERTIBLE NETWORKS

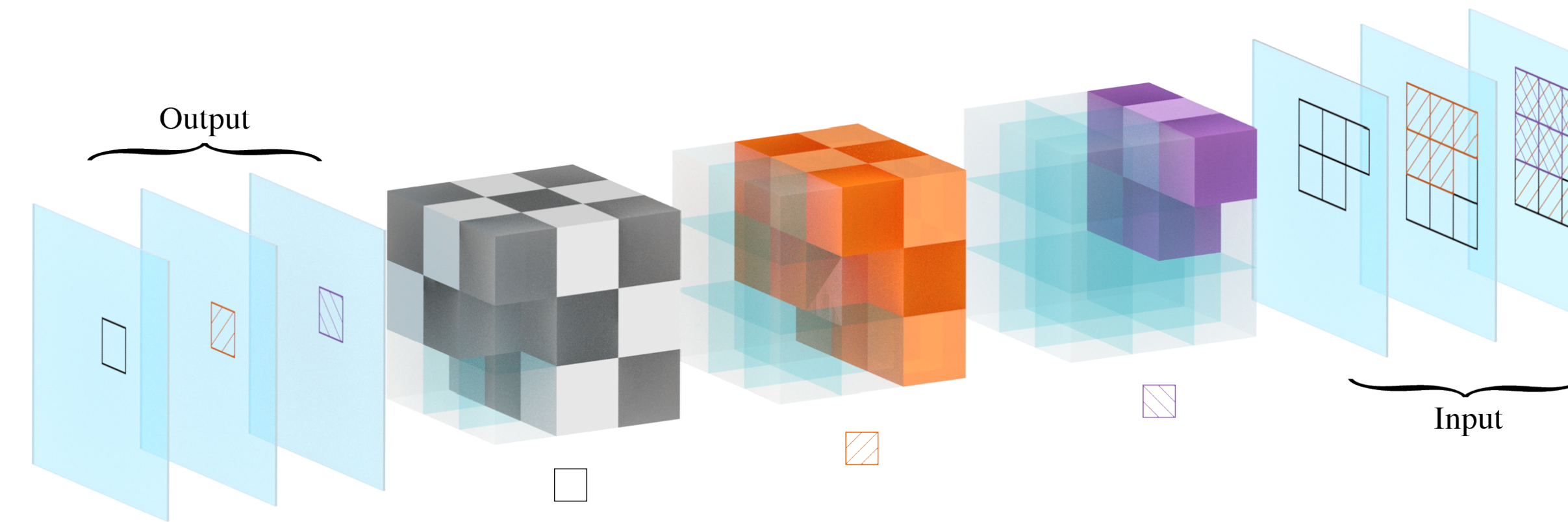


Figure 2: Illustration of a masked convolution with 3 filters and kernel size 3×3 .

Masked Invertible Layer (Mint Layer):

$$\mathbf{m}(\mathbf{x}) = \mathbf{t} \odot \mathbf{x} + \sum_{l=1}^K \mathbf{W}_i^3 h \left(\sum_{j=1}^K \mathbf{W}_{ij}^2 h(\mathbf{W}_j^1 \mathbf{x} + \mathbf{b}_j^1) + \mathbf{b}_{ij}^2 \right) + \mathbf{b}_i^3$$

- h : monotonic activation function where $h' \geq 0$.
- $\text{diag}(\mathbf{W}_i^3) \text{diag}(\mathbf{W}_{ij}^2) \text{diag}(\mathbf{W}_j^1) \geq 0, \forall 1 \leq i, j \leq K$
- $\{\mathbf{W}_i^1\}_{i=1}^K, \{\mathbf{W}_{ij}^2\}_{1 \leq i, j \leq K}$ and $\{\mathbf{W}_i^3\}_{i=1}^K$: masked convolutional layers.

Masked Invertible Network (MintNet):

- Paired Mint layers:** We always pair two Mint layers together: one with a lower triangular Jacobian and the other with an upper triangular Jacobian, so that the Jacobian of the paired layers is not triangular, and blind spots can be eliminated.
- Squeezing layers:** Subsampling is important for enlarging the receptive field of convolutions. We use an invertible “squeezing” operation to reshape the feature maps so that they have smaller resolution but more channels.

EXPERIMENTS

Classification:

Table 1: Classification accuracies

Model	MNIST (train)	MNIST (test)	CIFAR-10 (train)	CIFAR-10 (test)
ResNet	100%	99.6%	100%	92.6%
MintNet	100%	99.6%	100%	91.2%

Density estimation:

Table 2: Bits per dimension (bpd) results.

Method	MNIST	CIFAR-10	ImageNet 32×32
NICE [5]	4.36	4.48 [†]	-
MAF [22]	1.89	4.31	-
Real NVP [6]	1.06	3.49	4.28
Glow [15]	1.05	3.35	4.09
FFJORD [9]	0.99	3.40	-
i-ResNet [1]	1.06	3.45	-
MintNet (ours)	0.98	3.32	4.06

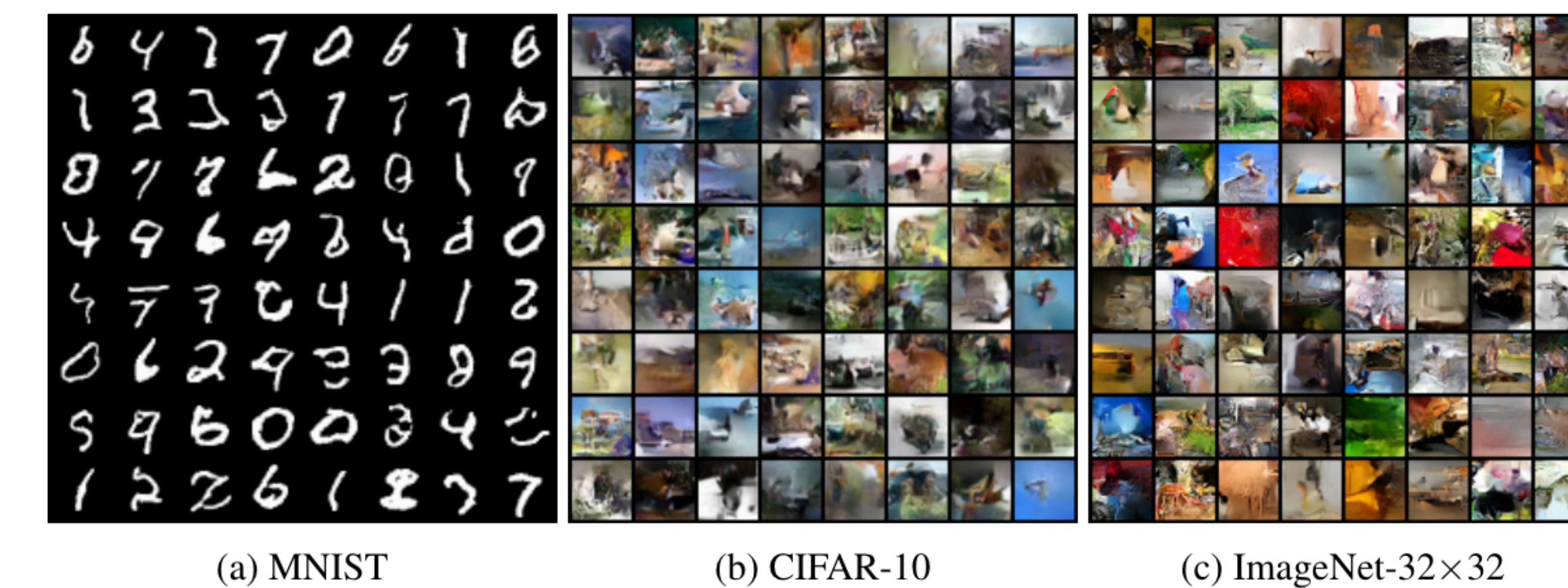


Figure 3: Uncurated samples.

Verification of invertibility:

We examine the performance of **Algorithm 1** by measuring the reconstruction error of MintNets.

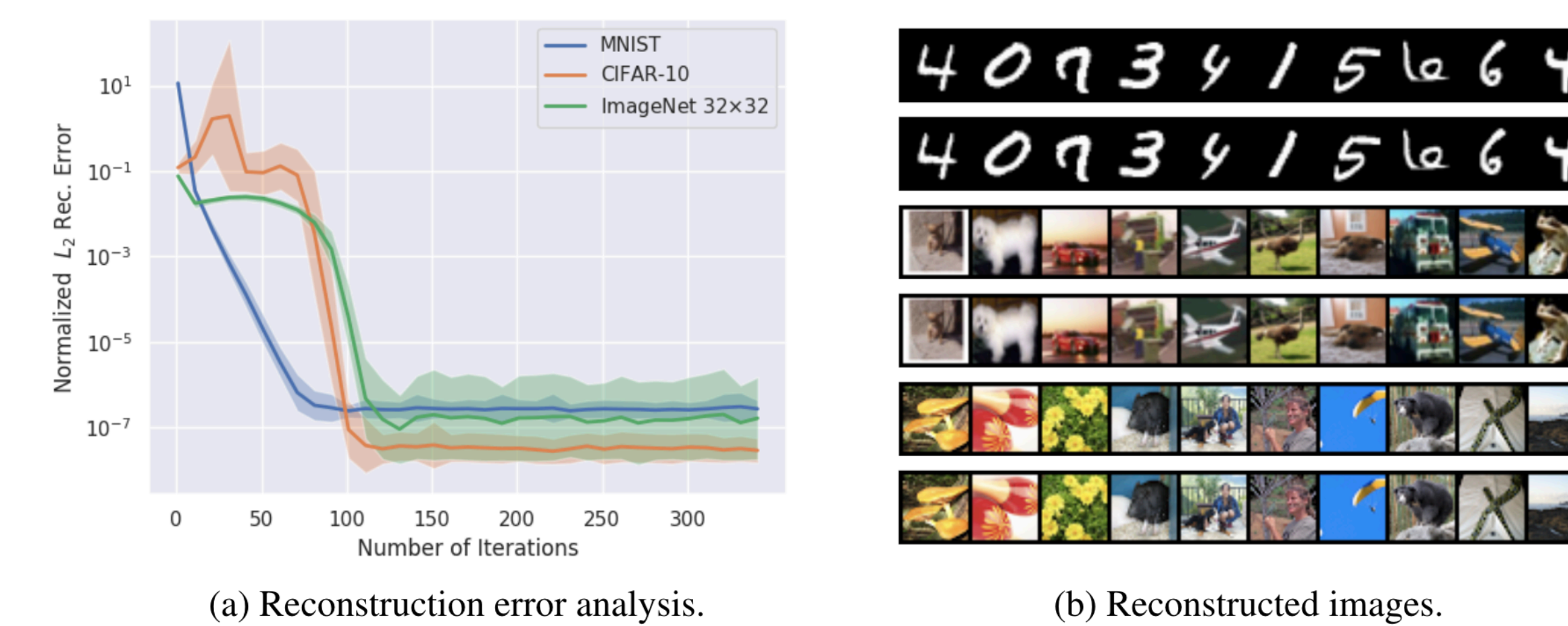


Figure 4: Accuracy analysis of **Algorithm 1**.

Interpolation of hidden representations:

Given four images $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4$ in the dataset, we interpolate over the feature space.



Figure 5: Interpolation of hidden representations.