# METAINV: OVERCOMING ITERATIVE AND DIRECT METHOD LIMITATIONS FOR INVERSE LEARNING

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Invertible neural networks (INNs) have gained significant traction in tasks requiring reliable bidirectional inferences, such as data encryption, scientific computing, and real-time control. However, iterative methods like i-ResNet face notable limitations, including instability on non-contractive mappings and failure in scenarios requiring strict one-to-one mappings. In contrast, analytical approaches like DipDNN guarantee invertibility but at the expense of performance, particularly in tasks demanding rich feature extraction (e.g., convolutional operations in complex image processing). This work presents a detailed analysis of the limitations in current invertible architectures, examining the trade-offs between iterative and analytical approaches. We identify key failure modes, particularly when handling information redundancy or strict bijections, and propose a meta-inverse framework that dynamically combines the advantages of both i-ResNet and DipDNN. Our framework adapts in real-time based on task-specific signals, ensuring both flexibility and guaranteed invertibility. Extensive experiments across diverse domains demonstrate that our hybrid approach outperforms existing methods in forward accuracy, inverse consistency, and computational efficiency. Our results highlight the utility of this meta-inverse strategy for critical applications where precision, stability, and adaptability are crucial.

027 028 029

030

025

026

004

010 011

012

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

In recent years, invertible neural networks (INNs) have become essential for a broad range of applications that require consistent bidirectional inference, such as data encryption, scientific computing, and real-time control systems Raissi et al. (2019); Devlin et al. (2019). In these settings, the ability to
precisely compute forward and inverse mappings is crucial for tasks like data recovery, state estimation, and adaptive control. Ensuring robust, accurate inverse computations underlies the reliability of these applications Levine et al. (2016); Lewis et al. (2012).

While iterative approaches such as i-ResNet Behrmann et al. (2019) offer flexibility and have been successful in complex, high-dimensional mappings (e.g., image processing and generative models), they come with significant drawbacks. i-ResNet depends on the Lipschitz condition for convergence, which is often difficult to guarantee in practice, especially in tasks requiring strict one-to-one mappings or non-redundant information. Furthermore, i-ResNet suffers from performance issues, particularly when the mappings are non-contractive, leading to unstable or inaccurate inverse estimations Zhang et al. (2020a). In some cases, it may even produce outputs that fall outside the valid range, further degrading its reliability in critical systems Kingma & Dhariwal (2018).

On the other hand, analytical methods like DipDNN (Yuan et al., 2024) ensure strict one-to-one invertibility at every layer through architectural constraints Dinh et al. (2014). These models are particularly well-suited for applications demanding high precision and exact inversions, such as realtime control and physics simulations. However, this architectural rigidity comes at a cost. DipDNN often underperforms in scenarios where flexible feature extraction—especially through convolutional layers—is necessary, such as in complex image processing tasks. The lack of convolutional layers, which are key for hierarchical feature learning, limits DipDNN's expressiveness and ability to handle high-dimensional, intricate data structures.

Given the complementary nature of iterative and analytical inverse methods, we argue that a single model may not perform optimally across all scenarios. Our extensive analysis of i-ResNet and

I = Reversible Data Hiding x + h(x) + h(

Figure 1: Contrast and Complement of i-ResNet's iterative inverse approximation and DipDNN's closed-form inverse. **MetaInv.** A switching algorithm to adapt model on diverse inverse problems.

DipDNN reveals that each method excels in different tasks but also suffers from critical shortcomings in others. Motivated by this observation, we propose a novel meta-inverse (MetaInv in Fig. 1) framework that dynamically switches between i-ResNet and DipDNN based on task-specific requirements. This approach leverages the strengths of both architectures—using i-ResNet for tasks requiring flexibility and feature extraction, while relying on DipDNN for tasks demanding strict, one-to-one invertibility. By incorporating a task-driven signal to select between models during inference, the framework adapts to the demands of different applications, ensuring robustness, precision, and computational efficiency.

Over meta-inverse framework provides a more generalized solution for a wide range of inverse learning tasks. The flexibility and adaptability of our approach allow for broader applicability, making it particularly useful in real-world scenarios that require both invertibility guarantees and high performance in feature-rich domains. In this paper, we present a thorough evaluation of our framework across various benchmarks, demonstrating that it significantly outperforms both i-ResNet and DipDNN individually. The results underscore the utility of combining iterative and analytical methods to tackle the diverse challenges faced in inverse learning.

086 087

088

054

055

057

060

066 067

068

069

#### 2 RELATED WORK

Invertible Neural Networks (INNs). Invertible neural networks (INNs) have emerged as essential tools for a variety of applications that require consistent bidirectional inferences, such as inverse 090 problems, density estimation, and generative modeling. One of the earliest models, NICE Dinh et al. 091 (2014), proposed a coupling layer design that ensures exact invertibility, followed by extensions like 092 RealNVP Dinh et al. (2016), which allowed more expressive feature transformation through affine coupling. Glow Kingma & Dhariwal (2018) further improved these architectures by incorporating 1x1 094 convolutional layers, making it more suitable for high-dimensional image tasks. Another prominent 095 line of research is iResNet Behrmann et al. (2019), which introduces iterative inversion based on 096 residual blocks, providing flexibility for complex mappings. However, these models have limitations in handling strict one-to-one mappings and suffer from numerical instabilities when dealing with 098 non-contractive mappings.

099 Analytical vs. Iterative Inversion Methods. In the context of inverse problems, there are two main 100 categories of methods: iterative and analytical. Iterative methods like iResNet Behrmann et al. (2019) 101 rely on residual connections and iterative approximation to compute the inverse. While this approach 102 is highly flexible and well-suited for complex tasks, it often struggles with stability and may produce 103 inaccurate inverse mappings, particularly in scenarios where strict one-to-one correspondences are 104 required. On the other hand, analytical inversion methods, such as DipDNN Dupont et al. (2019), 105 enforce strict bijection through architectural constraints. This guarantees invertibility but at the cost of flexibility, particularly in tasks requiring convolutional feature extraction. Analytical approaches 106 also suffer from reduced performance in tasks that require high-dimensional feature extraction, such 107 as image generation and classification.

108 Meta-Learning and Hybrid Architectures. Meta-learning has been explored as a way to dynam-109 ically adapt neural network architectures based on the task requirements. Recent works in online 110 control Christianson et al. (2023b) and hybrid models Blum & Burch (1997b) have demonstrated 111 the potential of combining model-based and learning-based approaches. This inspired our hybrid 112 meta-inverse framework, which leverages the strengths of both iterative and analytical methods. Specifically, our method dynamically switches between iResNet and DipDNN layers based on 113 task-specific cues, improving performance in diverse inverse learning tasks. This is similar in spirit 114 to works such as AdaNet Cortes et al. (2017) that adaptively grow neural architectures based on 115 performance metrics, but our approach focuses on invertibility and bidirectional consistency. 116

117 Applications in Real-Time Systems and Inverse Problems. The need for accurate and efficient 118 inverse mappings is especially pronounced in applications such as scientific computing Raissi et al. (2019), real-time control Levine et al. (2016), and robotics Ziebart et al. (2021). These domains 119 require models that can ensure robustness, precision, and real-time performance, particularly when 120 working with non-contractive or non-bijective mappings. While iterative methods like iResNet 121 provide flexibility, they often fail in applications with strict requirements for inverse consistency 122 and accuracy, such as fluid dynamics and physical system identification. Conversely, analytical 123 methods like DipDNN excel in these scenarios by providing exact inversions, although they lack 124 the flexibility required for more complex tasks like image generation or tasks that involve redundant 125 information. Our meta-inverse framework addresses these limitations by automatically selecting the 126 best method based on the problem requirements, leading to improved robustness and efficiency in 127 real-time applications.

128 129

130 131

132

140

141

142

147

148 149 150

151

152

153

154

#### 3 BACKGROUND: INVERSE PROBLEM AND INVERTIBLE MODELS

#### 3.1 PROBLEM STATEMENT AND DISTINCT DEMANDS

Inverse problems span numerous applications that involve recovering original variables from observed outputs This work focuses on inverse learning through invertible mapping recovery, which is for point estimates of images or physical system states. The problem is typically formulated as approximating a forward mapping  $f_{\theta} : \mathbb{R}^n \to \mathbb{R}^n$ , where  $y = f_{\theta}(x)$  is invertible. The goal is to find the relative inverse mapping  $g_{\vartheta}$  such that  $x = g_{\vartheta}(y) = f_{\theta}^{-1}(y)$ , ensuring consistency with the forward process. The demands for such two-way mapping rule recovery are distinct and varied even in one task. Thus, this work evaluates them using the following performance metrics.

• Forward Prediction Error (Fwd): Same as common one-way learning, it measures the ability of the model to predict y from x. The invertible model is trained by minimizing the forward prediction loss (any discriminative learning loss  $\ell_{\text{fwd}}$ ):

$$f^* = \operatorname*{argmin}_{ heta} \ell_{\mathrm{fwd}} \left( oldsymbol{y}, f_{ heta}(oldsymbol{x}) 
ight).$$

• Inverse Reconstruction Error (Inv): The reconstruction error evaluates the model's invertibility to recover the original inputs ( $\ell_{inv}$  is mean square error for point estimates),

$$\ell_{\mathrm{inv}}\left(oldsymbol{x},g_{artheta}\left(f_{ heta}(oldsymbol{x})
ight)
ight)$$
 .

- **Inverse Consistency (Inv-Consist**): It assesses the consistency between forward and inverse mappings by comparing the inverse predictions with the true labels *y*, rather than just the forward outputs:
  - $\ell_{\mathrm{fwd}}\left(oldsymbol{x},g_{oldsymbol{ heta}}(oldsymbol{y})
    ight).$

It is important to note the distinction between inverse accuracy and consistency. Inverse accuracy, or reconstruction error, indicates only the model's invertibility to compute the inverse, where analytical invertibility yields an error-free result, and numerical invertibility minimizes round-off errors. Consistency, on the other hand, evaluates the model's capability of two-way learning. It reveals the precision of global inversion, e.g., minimizing reconstruction error doesn't necessarily need a low forward approximation error, while the consistency error integrates errors from both forward and inverse processes. Most works involving approximate one-to-one mappings focus on inverse accuracy alone, such as the image classification and recovery (Behrmann et al., 2019).

## 162 3.2 REVIEW OF INVERSE LEARNING MODELS

Existing invertible neural networks (INNs) can be broadly categorized into two approaches: iterativenumerical approximations and analytical inverse solutions.

166

167 Iterative Numerical Methods: I-ResNet (Behrmann et al., 2019) and related models compute 168 inverse mappings iteratively. I-ResNet guarantees numerical stability in inversion by designing 169 residual blocks  $f_{i\text{-ResNet}}^{(k)}(z^k) = z^{k+1} = z^k + h^{(k)}(z^k)$ , where  $h^{(k)}$  is a nonlinear function that 170 satisfies  $\text{Lip}(h^{(k)}) < 1$  to ensure a contractive transformation. While this approach is computationally 171 flexible, the inverse  $z^k = (f^{(k)})^{-1}(z^{k+1})$  is obtained through iterative refinement, which can be 173 inefficient and prone to convergence issues, especially in scenarios where the Lipschitz condition 174 fails to hold.

174 175

Analytical Inverse Methods: These methods aim to provide closed-form inverses through algebraic constructions. For instance, NICE (Dinh et al., 2014) and Glow (Kingma & Dhariwal, 2018) achieve bijective transformations by using coupling layers that split input and output variables. The forward transformation is given by:  $z_1^{k+1} = az_1^k$ ,  $z_2^{k+1} = bz_2^k + t(z_1^k)$ , and the inverse is computed as:  $z_1^k = \frac{z_1^{k+1}}{a}$ ,  $z_2^k = \frac{z_2^{k+1} - t(z_1^k)}{b}$ . Here,  $a, b \neq 0$  are constants, and  $t(\cdot)$  is a nonlinear function (e.g., NN). While this ensures an analytical inverse, it imposes significant structural constraints that can limit expressiveness, especially in complex, high-dimensional tasks such as feature extraction and image modeling.

**Decomposed Invertible Pathway DNN (DipDNN):** DipDNN (Yuan et al., 2024) avoids strict partitioning by introducing a layer-wise decomposition with triangular weight matrices. The forward mapping in DipDNN is expressed as:

$$f_{\text{DinDNN}}^{(k)}(\boldsymbol{z}^k) = g_2^{(k)}(W_{\text{trij}}^{(k)}g_1^{(k)}(W_{\text{trij}}^{(k)}\boldsymbol{z}^k + b_1^{(k)}) + b_2^{(k)}).$$

The triangular structure of  $W_{tril}$  and  $W_{tri}$  ensures strict bijection, while the use of monotonic activation functions such as LeakyReLU or ELU preserves expressiveness. DipDNN offers a more flexible architecture that avoids the constrained coupling layers of NICE and Glow, enabling efficient inversions without sacrificing approximation power.

194

196

185

186

187 188 189

#### 4 INVERTIBILITY VS. APPROXIMATION LIMITATIONS OF I-RESNET

197 I-ResNet and its variants have achieved substantial success in tasks such as image classification and density estimation by addressing the challenge of invertibility. However, these 199 models struggle with the trade-off between expressive power 200 and the need for stable inversion, imposed by the Lipschitz 201 condition. This constraint forces each residual block to make 202 minimal adjustments to the input, limiting the model's ability 203 to approximate complex mappings that involve significant 204 changes or intricate dependencies (Zhang et al., 2020b). This 205 section provides theoretical analysis to explore these limita-206 tions in depth. 207

# 208 4.1 CONTRASTING INVERTIBILITY BETWEEN 209 I-RESNET AND DIPDNN



Figure 2: Fwd approximation capability on synthetic example.

Layer-Wise Contraction and Limited Scaling. The following theorem formalizes the limitation of i-ResNet in terms of its ability to approximate functions with high Lipschitz constants:

Theorem 1. For a K-layer i-ResNet, where each residual block  $h^{(k)}$  satisfies the Lipschitz constraint Lip $(h^{(k)}) < 1$ . Let  $f_T : \mathbb{R}^d \to \mathbb{R}^d$  be a multi-dimensional mapping that may be nonlinear, with a Lipschitz constant Lip $(f_T) = L_T > 1$ . Then, a K-layer i-ResNet cannot adequately approximate  $f_T$ if  $L_T > 2^K$ , and the relative approximation error is bounded by:  $\epsilon = \Omega \left(1 - \frac{2^K}{L_T}\right)$ . The proof are included in Appendix A.1 and A.2. This result indicates that as the complexity of the target function increases (i.e., as  $L_T$  grows), i-ResNet's ability to represent it deteriorates unless the number of layers is increased accordingly. For real-world applications, the number of layers required can be approximated as:

$$K_T \ge \frac{\log(\hat{K})}{\log(1 + \operatorname{Lip}(f_{i-\operatorname{ResNet}}))}$$

where  $\hat{K} = \max_i \frac{\max(y_i) - \min(y_i)}{\max(x_i) - \min(x_i)}$ . Detailed derivations of this estimate can be found in Appendix A.2. Empirical evidence suggests that adaptive layer estimation is effective when evaluated on small data batches.

**Inadequacy for Sign-Flipping.** While deeper architectures can mitigate the limitation of layer-wise contraction, i-ResNet faces an additional challenge in modeling mappings that involve sign-flipping in the inputs. The following theorem establishes that i-ResNet cannot approximate functions where input and output vectors have opposite signs. The corresponding proof is included in Appendix A.3.

Theorem 2. For a K-layer i-ResNet, where each residual block  $h^{(k)}$  satisfies the Lipschitz constraint Lip $(h^{(k)}) < 1$ , let  $f_T : \mathbb{R}^d \to \mathbb{R}^d$  be a multi-dimensional mapping such that  $\mathbf{x} \cdot f_T(\mathbf{x}) < 0$  (i.e., the signs of  $\mathbf{x}$  and  $f_T(\mathbf{x})$  are opposite in at least one dimension). We claim that i-ResNet  $f_{i-ResNet}$  cannot approximate such a mapping  $f_T$ .

236 Non-Strict One-to-One Mapping Approx-237 imation via Iterative Inverse. The layer-238 wise contraction of i-ResNet limits the degree 239 to which input space can be scaled or bent during forward transformations. However, this 240 approach to enforcing invertibility provides 241 flexibility in the network's architecture, e.g., 242 convolution layers within the residual blocks. 243 Specifically, the convolution operation is non-244 invertible due to its many-to-one mapping: 245  $y[i] = (x * w)[i] = \sum_{j=1}^{k} w[j] \cdot x[i-j],$ 246



Figure 3: (a) & (b) MNIST & CIFAR-10 classification accuracy; (c) Physical system modeling performance.

where w is the convolution filter or kernel. With padding, convolution aims to extract features from redundant information in the input, which is advantageous for tasks like classification involving high-dimensional images. In particular, convolution helps distill features and reduce dimensionality, improving performance (*e.g.*, Fig. 3 shows improved accuracy due to dimensionality reduction.

Nonetheless, there are many inverse problems where the redundancy of the information is minimal, 251 such as in fluid dynamics, sonar sensor assimilation, and power flow analysis. In these cases, either 252 the system states (spatial position, time) or physical properties (velocity, pressure, etc.) are loosely 253 dependent on one another. Importantly, physical systems typically have explicit forward functions, 254 and the ability to produce accurate point estimates for modeling and control is more critical than 255 generating qualitative results such as images. In fluid flow modeling, for instance, Navier-Stokes 256 equations govern the system, and any loss of information leads to ambiguous results or reconstruction 257 errors. Under such scenarios, as Fig. 3 (c) indicates, i-ResNet's iterative approximation of non-strict 258 one-to-one mappings may result in inaccuracies and lack of exactness in the inverse computations.

259

220 221 222

223 224

225

226 227

228

229

230

Layer-Wise Bijective Transformation for Analytical Inverse. In contrast to the iterative nature of i-ResNet, DipDNN enforces invertibility through layer-wise one-to-one transformations. The proposition to guarantee bijective layer-wise transformation is in Appendix A.4. Unlike i-ResNet's Lipschitz constraints, which limit flexibility, DipDNN imposes no limitations on the network's ability to scale features or flip the signs of inputs.

265

Equivalence & Conflict Between i-ResNet and DipDNN. While the strict bijective transforma tions of DipDNN provide guaranteed invertibility, this structure can be overly restrictive for some
 inverse problems where approximate inverses suffice. i-ResNet, with its flexibility (e.g., convolution
 layers), may outperform DipDNN in these cases. This raises an important question: can DipDNN
 achieve similar performance for these tasks?

273

274 275

276

277 278

279

281

282

300

304

305

306

307

308

310

311

270 To address this, we aim to show an equivalence between convolution-based residual networks and 271 DipDNN based on fully connected layers, for which the proof is in Appendix A.5. 272

**Theorem 3.** Consider a K-layer convolution-based i-ResNet. Let each convolution layer  $h_i^{(k)}(z^k) =$  $g(W_{conv}^{(k)} * \boldsymbol{z}^k + b)$ , where  $W_{conv} \in \mathbb{R}^{k \times k}$ . Then:

- 1) Each residual block  $h_i^{(k)}$  can be represented in an equivalent form in DipDNN, satisfying Proposition 1.
- 2) The residual connection  $f_{i\text{-ResNet}}^{(k)}(\boldsymbol{z}^k) = \boldsymbol{z}^k + h^{(k)}(\boldsymbol{z}^k)$  inherently violates the strict invert*ibility of DipDNN.*

#### META-INVERSE ALGORITHM FOR COMPLEMENTARY SELECTION 4.2

283 From Sections 4 and 4.1, we observe the contrasting strengths and weaknesses of i-ResNet and 284 DipDNN. i-ResNet excels in providing flexible architectures, such as non-bijective convolutional feature extraction, but struggles to maintain strict one-to-one mappings and compensates poorly when faced with contraction limitations. On the other hand, DipDNN guarantees analytical invertibility 286 through its structured, layer-wise design, but this strictness can be overly limiting for inverse problems 287 that exhibit redundancy in the data. These differences highlight the need for a hybrid approach capable 288 of dynamically adapting to task-specific characteristics, e.g., complexity, redundancy, and precision 289 requirements of inverse problems. 290

To address this, we propose the Meta-Inverse 291 (MetaInv) algorithm, which dynamically switches be-292 tween i-ResNet and DipDNN based on the specific 293 characteristics of the task at hand. The MetaInv algorithm in Fig. 4 is inspired by learning-augmented 295 switching methods in online control (Christianson 296 et al., 2023a; Blum & Burch, 1997a), where an agent al-297 ternates between a learning-based approach with high 298 expressive power (i-ResNet) and a model-based ap-299 proach with lower risk (DipDNN).



Figure 4: Meta-Inverse Algorithm: Switching between i-ResNet and DipDNN based on task characteristics.

- Specifically, as outlined in Section 3.1, the MetaInv 301 algorithm is designed to balance four key requirements of inverse learning: 302 303
  - (a) Forward expressive power: Modeling complex mappings correctly with high fidelity.
  - (b) Tight error bounds on the inverse: Minimizing inverse reconstruction error, ideally achieving values close to zero (theoretically) or minimizing numerical round-off errors.
  - (c) Consistency of bi-directional mappings: Ensuring that both forward and inverse mappings are consistent when true labels are used for inverse predictions.
    - (d) Computational efficiency: Optimizing resource usage to ensure that the model is both computationally efficient and scalable to large datasets.

**Trust-Weighted Switching Mechanism:** The core of the MetaInv algorithm (Algorithm 2) is a 312 trust-weighted switching mechanism that integrates task-specific performance metrics (Fwd, Inv, 313 Inv-Consist) defined in Sec. 3.1 with computational cost into a single evaluation function:  $V_{\text{model}} =$ 314  $J_{\text{model, total}} + \lambda C_{\text{model}}$ , where  $J_{\text{model, total}}$  is the weighted combination of performance metrics,  $\lambda$  is the 315 weight for computational cost, and  $C_{\text{model}}$  includes time complexity and model size. 316

To ensure stability and prevent trivial convergence, the trust-weight parameter  $\beta$  is updated iteratively 317 based on the performance difference between i-ResNet and DipDNN:  $\beta_{t+1} = \beta_t + \eta_t (V_{i-\text{ResNet}} - \beta_t)$ 318  $V_{\text{DipDNN}}$ ), where  $\eta_t$  is the learning rate (empirically set to 0.01). Additionally,  $J_{\text{threshold}}$  sets an 319 acceptable range of performance, preventing the algorithm from stagnating when inverse learning is 320 sensitive to noise or data uncertainties. This mechanism allows the algorithm to smoothly adapt to 321 the optimal model for an arbitrary task. 322

- The evaluation of task-specific metrics is not fixed and evolves dynamically during learning dynami-323 cally switches between the two architectures, optimizing both forward and inverse computation. This
  - 6

4	Algorithm 1 Meta-Inverse Algorithm Between i-ResNet and DipDNN
:5 :6	<b>Input:</b> $\{x_i, y_i\}_{i=1}^N$ (data), $J_{\text{threshold}}$ (performance threshold), $\lambda$ (weights for performance and cost),
7	T (iterations)
8	<b>Output:</b> Optimal selection of i-ResNet or DipDNN for computations
0	Initialize $\beta_0$ , $\lambda$
.9	for $t = 1, \dots, T$ do
U	Forward and Inverse Computation
1	Compute forward, inverse (reconstruction and ablation) outputs: $y_{\text{model}}$ , $\hat{x}_{\text{model}}$ , $\hat{x}'_{\text{model}}$ for
2	$model \in \{i-ResNet, DipDNN\},\$
3	Metric Evaluation
4	Calculate Fwd Acc, Inv Acc, and Inv Consist losses and compute combined scores:
5	$J_{\text{model, total}} = \sum_{k=1}^{3} \lambda_k J_{\text{model, k}}$ , where $J_{\text{model, k}}$ represents the respective metrics for $k = 1, 2, 3$ .
6	Model Selection
7	if $J_{\text{DipDNN, total}} < J_{\text{i-ResNet, total}}$ and $J_{\text{DipDNN, total}} < J_{\text{threshold}}$ then
8	Use DipDNN: $y = y_{\text{DipDNN}}, x_{\text{reconstructed}} = x_{\text{DipDNN}}$
9	else
0	Use i-ResNet: $y = y_{i-ResNet}$ , $x_{reconstructed} = x_{i-ResNet}$
1	end if
0	Trust-Weight Update
2	Compute computational cost $C_{\text{model}}$ for both models: $C_{\text{model}} = \alpha_1 T_{\text{model}} + \alpha_2 M_{\text{model}} + \alpha_3 S_{\text{model}}$
3	Compute total evaluation cost $V_{\text{model}}$ for each model: $V_{\text{model}} = J_{\text{model, total}} + \lambda C_{\text{model}}$
4	Update trust-weight $\beta_{t+1} = \beta_t + \eta_t (V_{i-\text{ResNet}} - V_{\text{DipDNN}})$
5	end for

approach corrects the limitations of prior evaluations that were either incomplete or misaligned with the actual data characteristics and task-specific requirements. Details of the MetaInv algorithm's implementation and evaluation can be found in the Appendix.

355 356

350

347 348 349

#### 5 NUMERICAL EXPERIMENTS

357 Throughout experiments on diverse inverse problems, we aim to answer the following questions: 1) 358 For non-contractive cases (scaling & sign-flipping), can DipDNN outperform i-ResNet on bidirectional mappings? 2) Iterative inverse approximation vs. analytical inverse on non-strict bijective cases 359 of different information redundancy, which wins out? 3) Can MetaInv switch to the optimal model 360 for arbitrary Inverse problems? To answer these questions, we perform our experiments on different 361 inverse problems: synthetic examples, image classification (Fwd) and reconstruction (Inv) from last 362 features (Behrmann et al., 2019), image transformation with inverted colors for reversible data hiding 363 (Zhang et al., 2024), power flow modeling (Fwd) and state estimation (Inv), and Navier-Stokes flow 364 dynamics modeling (Fwd) and initial condition inference (Inv) (Langtangen & Logg, 2017). Each task represents important application needs in inverse problems. 366

For comparison, we test 1) **i-ResNet**, which is representative of iterative inverse approximation; 367 2) **ResNet**, a baseline to reveal the standard approximation capability compared with contractive 368 i-ResNet; 3) NICE, a baseline analytical inverse model with strict architecture constraint; 4) DipDNN, 369 an improved analytical inverse with layer-wise bijective transformation. Finally, we compare them 370 with the MetaInv, the proposed switching algorithm to adaptively select between i-ResNet and 371 DipDNN. Note that we select the most representative invertible NN structures, whose variants have 372 not shown significant differences for the targeted broad applications in both image-related and 373 physical system tasks. For example, although Glow (based on NICE) is not directly tested, we 374 implement its invertible ActNorm and  $1 \times 1$  convolution in experiments. For all the testing scenarios, 375 all the inverse learning models contain similar numbers of parameters in each comparison for a fair evaluation. For each inverse learning method, we equally evaluate the performance using the error 376 metrics of forward prediction, inverse reconstruction, and inverse consistency, as defined in Sec. 3.1. 377 Specially for image tasks, we show qualitative results for intuitive visualization.

#### 5.1 INVERTIBILITY VS. EXPRESSIVE POWER ON NON-CONTRACTIVE CASES

Based on the theoretical analysis in Sec. 4, i-ResNet imposes limitations in expressive power. We start with synthetic datasets to quantify the limitations. Further, we conduct sensitivity analysis and visualize qualitative results to understand the trade-off between invertibility and approximation capability.

384 385

378

379 380

381

382

386 387

388

389

390

405

414 415 Sensitivity Analysis of Invertible NN Depth. To test the number of layers needed for i-ResNet, we evaluate the performance of synthetic data (Fig. 2) and power flow cases (Fig. 5), which exhibit non-contractive mappings. It's obvious that while DipDNN approaches a similar fwd error as ResNet, i-ResNet has a much higher error. It requires a deeper model to reach the same performance as others.

**Color Inversion Mixed with Complex Transformation.** 391

The invertibility enforced by contraction imposes the inabil-392 ity to approximate xf(x) < 0. Except for the synthetic data 393 with the flipped sign of inputs and outputs (Fig. 10), we 394 generalize such a property to more complex image examples. 395 It is a reversible data-hiding task that mixes color inversion 396 (e.g., white to black and black to white) and class transfor-397 mation (e.g., one item to another) to cover the privacy infor-398 mation of original images Fig. 6. The data-hiding process needs to be lossless to recover back to the original data. 399

In the synthetic example with flipped signs, Fig. 10 in the 400 Appendix shows that the error of i-ResNet cannot decrease 401 no matter how many more layers are used. Other invertible 402 NNs successfully capture this relationship but NICE cannot Figure 5: Fwd approximation capabil-403 maintain the inverse performance on both reconstruction ity on power flow example. 404 (middle column) and independent prediction (last column).



In a similar context, we construct the data-hiding case in images (Fig. 6) using images from MNIST, 406 FashionMNIST, and CIFAR-10. Regular images are used as inputs; we create outputs by combining 407 half of the color-converted MNIST images with half of the FashionMNIST images and, alternatively, 408 a quarter of color-converted MNIST images with three-quarters of FashionMNIST images. When 409 gradually mixing image class transformation and color inversion (a portion of  $\{0, 1/2, 1/4\}$ ), i-410 ResNet is hard to find the correct mapping. Specifically, the 3rd and 4th rows in Fig. 6 indicate 411 that i-ResNet merely learns. Analytical invertible NNs like NICE and DipDNN not only have no 412 restrictions but also learn the nonlinear transformation well. Moreover, DipDNN has superior inverse 413 consistency for balanced performances of bi-directional mappings.



Figure 6: Performance of different models on color inversion mixed with complex transformation.

430 431

## 432 433 434 5.2 APPROXIMATED INVERSE VS. CLOSED-FORM INVERSE ON DIFFERENT INFORMATION REDUNDANCY 434

In Sec. 4.1, we analyze that the performances of invertible NNs depend on the data and problem need. Thus, we test on two representative physical and engineering systems other than the previous image-related tasks. Unlike the image classification case, the data of power system measurements and fluid dynamics observations have limited information redundancy as they are clearly defined system states or conditions.



Figure 7: Performance comparison of different methods on the power flow case. Top-left: forward voltage predictions (Fwd); Bottom plots: inverse load recovery (Inv and Inv-Consist); Top-right Table: numerical error metrics. I-ResNet achieves tightly bounded errors but predicts deviate from the ground truth. DipDNN captures system dynamics more effectively, which MetaInv selected to balance predictive accuracy and inverse alignment.

As for the power flow case, the data is generated from the simulation of grid power flow equations, which govern the system-wide power injections based on voltage phasors and the system's topology and parameters. The power flow dataset is built on the load inputs to estimate the voltage at different buses in an 8-bus system for forward mapping. In this case, inverse learning is essential to recover the underlying system physics and estimate load conditions for better interpretability and consistency in power flow analysis.

The power system experiment shows how different methods predict voltage (Fwd) and load (Inv) under limited measurements Fig. 7. DipDNN and MetalInv outperform i-ResNet using the iterative inverse. Although their errors are similar in the table at the top-right corner, DipDNN fits individual points much better. This is because i-ResNet has the contractive property enforced by Lip < 1: even though the point estimates are bad, the average error is bounded tightly. DipDNN is selected by MetaInv for the power flow case. The voltage plot (Fwd) reveals that I-ResNet, despite contracting the error, produces a pattern that does not align with the ground truth. Conversely, DipDNN achieves better consistency and captures the essential dynamics of the system, especially in the inverse problem (Inv-Consist). This demonstrates how critical it is for the method to maintain a balance between predictive power and inverse accuracy. Moreover, convolutional layers, while useful for forward problems, might introduce excessive contraction in systems with low redundancy, leading to distortions in the inverse process. 

As for the fluid dynamics case, the data is generated from the simulation of a chemical reaction
 between species A and B, which form species C, while considering advection and diffusion in a
 domain governed by a coupled system of partial differential equations. The dataset is built on the



Figure 8: Forward and inverse performance of different models on Navier-Stokes flow dynamics.



Figure 9: Performance of convolutional models on Navier-Stokes flow dynamics.

concentrations of these species, along with the velocity field, to capture the spatio-temporal evolution of the chemical reaction and the dynamics of fluid flow. In this case, inverse learning is needed to recover system dynamics and parameters from observed concentration data for interpretability and consistency. With dense observation data, Fig. 8 depicts the surface plots for forward (Fwd) and inverse consistency (Inv-Consist) of the velocity fields (wx and uA) for multiple models. DipDNN shows a superior fit for both forward and inverse consistency, especially in cases where the information redundancy is higher, and the system is more over-determined. For i-ResNet, although it attempts to approximate the inverse, its surface predictions deviate significantly, reflecting its struggle with preserving the intricate dynamics of the fluid flow. 

The error plots comparing prediction and truth further reinforce this. DipDNN maintains a low mean
 squared error (MSE) across both forward and inverse mappings, while i-ResNet struggles with a high
 inverse consistency error. This discrepancy highlights the limitation of iterative inverse approximation,
 especially when the system's redundant information is minimal. In this case, DipDNN's layer-wise
 bijective transformations enable a more robust handling of both the advection and diffusion terms in
 the PDEs, capturing the fluid dynamics more precisely.

Figure 9 provides insights into the impact of convolutional layers when information redundancy
is low. As seen in the forward and inverse consistency surface plots, both ResNet and I-ResNet
produce distorted surfaces. Convolutional layers introduce locality in modeling, which is beneficial
with high information redundancy. However, in systems like Navier-Stokes dynamics with limited
information, the imposed locality leads to contracted errors and incorrect surfaces. The experiments
show that DipDNN, which avoids heavy convolutional operations, performs significantly better in
such low-redundancy cases.

### <sup>540</sup> 6 CONCLUSION

541 542

In this paper, we analyze the complementary strengths and inherent limitations of typical invertible 543 models DipDNN and i-ResNet, highlighting the challenges of applying either model individually 544 across a wide range of inverse problems. DipDNN ensures strict analytical inversion with precise one-to-one mapping but lacks the flexibility to incorporate convolutional layers for feature extraction, 546 whereas i-ResNet effectively integrates convolution with residual connections, enhancing feature extraction and scalability; however, its Lipschitz constraints limit forward approximation, compro-547 548 mising inverse accuracy in non-contractive scenarios. Thus, we propose a meta-inverse algorithm to leverage their respective strengths on specific inverse problems. MetaInv balances three critical 549 metrics-forward accuracy, inverse consistency, and inverse accuracy-addressing the shortcomings 550 of prior work that focused on a subset of these metrics. By coupling i-ResNet and DipDNN effectively, 551 our solution ensures a more comprehensive evaluation and a fairer comparison for diverse inverse 552 problems.

553 554 555

556

558

559

562

569

570 571

572

573

574

577

580

581

582

585

586

587 588

589

590

#### References

- Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, pp. 573–582, 2019.
- Avrim Blum and Carl Burch. On-line learning and the metrical task system problem. In *Proceedings of the tenth annual conference on Computational learning theory*, pp. 45–53, 1997a.
- Avrim Blum and Carl Burch. On-line learning and dynamic model selection. In *Information Theory*, 1997. *Proceedings.*, 1997 IEEE International Symposium on, pp. 318, 1997b.
- Nicolas Christianson, Junxuan Shen, and Adam Wierman. Optimal robustness-consistency tradeoffs
   for learning-augmented metrical task systems. In *International Conference on Artificial Intelligence and Statistics*, pp. 9377–9399. PMLR, 2023a.
  - Paul Christianson, Josip Djolonga, Andreas Georghiou, and Vasileios Tzoumas. Learning-augmented switching algorithms for online control. *arXiv preprint arXiv:2301.00450*, 2023b.
  - Corinna Cortes, Javier Gonzalvo, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Adanet: Adaptive structural learning of artificial neural networks. In *International Conference on Machine Learning*, pp. 874–883, 2017.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
   bidirectional transformers for language understanding. In *Proceedings of the NAACL-HLT*, 2019.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components
   estimation. *arXiv preprint arXiv:1410.8516*, 2014.
  - Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In *International Conference on Learning Representations*, 2016.
- Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. In *Advances in Neural Information Processing Systems*, pp. 3140–3150, 2019.
  - Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10215–10224, 2018.
  - Hans Petter Langtangen and Anders Logg. Solving PDEs in python: the FEniCS tutorial I. Springer Nature, 2017.
- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.
- 59

Frank L Lewis, Draguna Vrabie, and Vassilis L Syrmos. Optimal control. John Wiley & Sons, 2012.

- Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A
   deep learning framework for solving forward and inverse problems involving nonlinear partial
   differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- Jingyi Yuan, Yang Weng, and Erik Blasch. Dipdnn: Preserving inverse consistency and approximation
   efficiency for invertible learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4071–4082, 2024.
  - Cheng Zhang, Bo Ou, Fei Peng, Yao Zhao, and Keqin Li. A survey on reversible data hiding for uncompressed images. ACM Comput. Surv., 56(7), April 2024. ISSN 0360-0300. doi: 10.1145/3645105. URL https://doi.org/10.1145/3645105.
  - Guoxia Zhang, Jing Hu, and Shiwen Yu. Efficient and robust training of deep neural networks with reduced numerical precision. In *International Conference on Learning Representations*, 2020a.
  - Han Zhang, Xi Gao, Jacob Unterman, and Tom Arodz. Approximation capabilities of neural odes and invertible residual networks. In *International Conference on Machine Learning*, pp. 11086–11095. PMLR, 2020b.
    - Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Inverse reinforcement learning: A unified survey. *Proceedings of Machine Learning Research*, 2021.

#### A APPENDIX

#### A.1 PROOF OF THEOREM 1

*Proof.* Given that  $Lip(h^{(k)}) < 1$ , the Lipschitz constant of each residual layer satisfies  $Lip(f^{(k)}) < 2$ . Hence, for a K-layer i-ResNet, the overall Lipschitz constant of the network is bounded by:

$$\operatorname{Lip}(f) = \prod_{k=1}^{K} (1 + \operatorname{Lip}(h^{(k)})) < 2^{K}.$$

If the target function  $f_T$  has a Lipschitz constant  $L_T > 2^K$ , the i-ResNet is unable to approximate it accurately since its representation capacity is limited by  $2^K$ . The relative approximation error across  $\epsilon$  dimensions is thus bounded by:

$$\epsilon = rac{\|f(oldsymbol{x}) - f_T(oldsymbol{x})\|}{\|f_T(oldsymbol{x})\|} = \Omega\left(1 - rac{2^K}{L_T}
ight).$$

#### A.2 ESTIMATING THE NUMBER OF LAYERS NEEDED FOR I-RESNET

To approximate a target mapping  $f_T$  with  $f_{i-\text{ResNet}}$ , the effective Lipschitz constant  $\hat{K}$  can be empirically estimated from the data as:

$$\hat{K} = \max_{i} V_i, \quad V_i = \frac{\max(y_i) - \min(y_i)}{\max(x_i) - \min(x_i)}, \quad i = 1, \dots, d.$$

For an i-ResNet with K layers, the effective Lipschitz constant can be approximated as:

 $\operatorname{Lip}(f_{i-\operatorname{ResNet}}) \approx (1 + \operatorname{Lip}(h^{(k)}))^{K}.$ 

<sup>641</sup> Note that Lip $(h^{(k)})$  can be adjusted close to a number < 1 in experiments and we use it for estimation. To approximate  $f_T(x)$ , the number of layers  $K_T$  necessarily needed is estimated by:

$$K_T \ge \frac{\log(K)}{\log(1 + \operatorname{Lip}(h^{(k)}))}.$$

647 If  $\hat{K} \gg 1$ , then  $K_T$  has to be sufficiently large for forward accuracy, compared to ResNet (without constraints) and DipDNN.

#### A.3 SUPPLEMENTARY PROOF OF FLIP-SIGN MAPPING FOR I-RESNET IN THEOREM 2

*Proof.* Consider a target mapping  $f_T(x)$  where the dot product  $x \cdot f_T(x) < 0$ , indicating that the inputs and outputs have opposite signs in at least one dimension. To demonstrate that i-ResNet cannot model this mapping, we show that for any  $x, x \cdot f_{i-\text{ResNet}}(x) \ge 0$  holds. 

From (Zhang et al., 2020b), consider two inputs  $x_1^0$  and  $x_2^0 = x_1^0 + \delta_0$ . After one residual block, the outputs are  $x_1^1 = x_1^0 + h^{(1)}(x_1^0)$  and  $x_2^1 = x_1^0 + \delta_0 + h^{(1)}(x_1^0 + \delta_0)$ . The difference between  $x_1^1$  and  $x_2^1$  is given by:

$$\delta_1 = \boldsymbol{x}_2^1 - \boldsymbol{x}_1^1 = \delta_0 + \left( h^{(1)}(\boldsymbol{x}_1^0 + \delta_0) - h^{(1)}(\boldsymbol{x}_1^0) \right).$$

Since  $Lip(h^{(1)}) < 1$ , it follows that: 

implying that the transformation preserves the sign relationships between  $x_1^0$  and  $x_2^0$ . Applying this logic iteratively to deeper layers shows that for any input x, the residual layers maintain  $\boldsymbol{x} \cdot f_{\text{i-ResNet}}(\boldsymbol{x}) \geq 0.$ 

 $|\delta_1| < |\delta_0|,$ 

Thus, i-ResNet cannot model mappings where  $x \cdot f_T(x) < 0$ . The relative approximation error for such mappings is quantified by:

$$\epsilon = \inf_{f_{i\text{-ResNet}}} \frac{\|f_T(\boldsymbol{x}) - f_{i\text{-ResNet}}(\boldsymbol{x})\|}{\|f_T(\boldsymbol{x})\|} \ge \Omega(1),$$

where the error is dominated by the components with flipped signs.

#### A.4 GUARANTEE INVERTIBILITY OF DIPDNN

**Proposition 1.** The neural network model  $f_{DipDNN} : \mathbb{R}^n \to \mathbb{R}^n$ , defined as  $f_{DipDNN} = f_{DipDNN}^{(1)} \circ \cdots \circ$  $f_{DipDNN}^{(K)}$ , is invertible if the weight matrices  $W_{tril}^k$  and  $W_{triu}^k$ , for  $k \in [1, K]$ , are lower and upper triangular matrices with non-zero diagonal elements. Each block  $f_{DipDNN}^{(k)}(\boldsymbol{z}^k)$  is given by: 

$$f_{DipDNN}^{(k)}(\boldsymbol{z}^{(k)}) = g_2^k(W_{tril}^k g_1^k(W_{triu}^k \boldsymbol{z}^k + b_1^k) + b_2^k)$$

where  $W_{tril}^k, W_{triu}^k \in \mathbb{R}^{n \times n}$ , and the non-zero diagonal elements of  $W_{tril}^k$  and  $W_{triu}^k$  ensure invertibility via the triangular structure. The activation functions  $g_1^k$  and  $g_2^k$  are strictly monotonic, ensuring bijective transformations for each layer.

The inverse of the DipDNN transformation can be computed in closed form as follows:

$$f_{\text{DipDNN}}^{(k)^{-1}}(\boldsymbol{z}^{k+1}) = \left(W_{\text{triu}}^k\right)^{-1} \left(g_1^k\right)^{-1} \left(\left(W_{\text{tril}}^k\right)^{-1} \left(\left(g_2^k\right)^{-1} (\boldsymbol{z}^{k+1} - b_2^k)\right) - b_1^k\right)$$

A.5 PROOF OF THEOREM 3

*Proof.* For 1), we rewrite the convolution operation as matrix multiplication, with  $T(W_{conv}^{(k)})$  representing the convolution kernel as a Toeplitz matrix:

$$h_i^{(k)}(\boldsymbol{z}^k) = g(W_{\text{conv}}^{(k)} * \boldsymbol{z}^k + b) = g(T(W_{\text{conv}}^{(k)})\boldsymbol{z}^k + b) = g(L_{\text{conv}}^{(k)}U_{\text{conv}}^{(k)}\boldsymbol{z}^k + b)$$

Applying LU decomposition to the Toeplitz matrix, we obtain lower and upper triangular matrices, ensuring invertibility. Therefore, in DipDNN, we assign  $W_{\text{tril}}^k = L_{\text{conv}}^{(k)}, W_{\text{triu}}^k = U_{\text{conv}}^{(k)}$ , and  $g_1^k$  as identity, resulting in: 

$$f_{\rm DipDNN}^{(k)}(\boldsymbol{z}^k) = g_2^k(W_{\rm tril}^k g_1^k(W_{\rm triu}^k \boldsymbol{z}^k + b_1^k) + b_2^k) = g(L_{\rm conv}^{(k)} U_{\rm conv}^{(k)} \boldsymbol{z}^k) = h_i^{(k)}(\boldsymbol{z}^k).$$

For 2), the residual connection  $f_{i\text{-ResNet}}^{(k)}(\boldsymbol{z}^k) = \boldsymbol{z}^k + h^{(k)}(\boldsymbol{z}^k)$  introduces an additive identity term that violates DipDNN's strict triangular structure. The addition imposes ambiguity, as the nonlinear dependency between  $h_i^{(k)}(z)$  and z precludes a closed-form inverse solution. Thus, DipDNN's architecture conflicts with i-ResNet's residual addition, as DipDNN requires strict bijective mappings. 

## A.6 DETAILS OF META-INVERSE ALGORITHM IMPLEMENTATION

We take inspiration from trust-region policy optimization and learning-augmented switching algorithms for selecting control agents. MetaInv aims to dynamically select between i-ResNet (analogous to a learning-based, higher-expressive-power agent) and DipDNN (analogous to a model-based, lower-risk agent) similarly so that the overall system's performance is optimized with stability guarantees.

Algorithm 2 Meta-Inverse Algorithm Between i-ResNet and DipDNN	
<b>Input:</b> $\{x_i, y_i\}_{i=1}^N$ (data), $J_{\text{threshold}}$ (performance threshold), $\lambda$ (weights for performance and co	ost),
T (iterations)	
<b>Output:</b> Optimal selection of 1-ResNet or DipDNN for computations	
Initialize $\beta_0, \lambda$	
Ior $t = 1, \dots, I$ do Econyand and Inverse Computation	
Forward and inverse (reconstruction and shlation) outputs: $a_1 \dots \hat{m}'$	for
$y_{model}, x_{model}, x_{model}$	101
Metric Evaluation	
Calculate Fwd Acc, Inv Acc, and Inv Consist losses and compute combined sco	res:
$J_{\text{model, total}} = \sum_{k=1}^{\infty} \lambda_k J_{\text{model, k}}$ , where $J_{\text{model, k}}$ represents the respective metrics for $k = 1, 2$ . Model Selection	3.
if $J_{\text{DipDNN, total}} < J_{\text{i-ResNet, total}}$ and $J_{\text{DipDNN, total}} < J_{\text{threshold}}$ then	
Use DipDNN: $y = y_{\text{DipDNN}}, x_{\text{reconstructed}} = x_{\text{DipDNN}}$	
else	
Use i-ResNet: $y = y_{i-\text{ResNet}}, x_{ ext{reconstructed}} = x_{i-\text{ResNet}}$	
end if	
Trust-Weight Update	
Compute computational cost $C_{\text{model}}$ for both models: $C_{\text{model}} = \alpha_1 T_{\text{model}} + \alpha_2 M_{\text{model}} + \alpha_3 S_{\text{model}}$	odel
Compute total evaluation cost $V_{\text{model}}$ for each model: $V_{\text{model}} = J_{\text{model, total}} + \lambda C_{\text{model}}$ Update trust-weight $\beta_{t+1} = \beta_t + \eta_t (V_{\text{i-ResNet}} - V_{\text{DipDNN}})$ end for	
In the following, we summarize the setups of parameters in implementation.	
• $J_{\text{threshold}}$ indicates the acceptable value range of the inverse learning performance in cert tasks. The value is determined by cross-validation to prevent being trapped in the lo optimum when training invertible models. It is used as we observe both models r experience trivial convergences in training.	ain cal 1ay
• Weighting factors $\lambda_i$ , $I = 1, 2, 3$ are the most important to the switching results and task-specific to combine the three metrics of inverse learning. Without prior knowledge the datasets, we use equal weights $(\lambda_1 = \lambda_2 = \lambda_3)$	are of
We set the number of iterations for emitting $T = T_2 = T_3$ .	
• we set the number of iterations for switching $I = 500$ to allow the model to fully adap	1.
• The weights $\alpha_1, \alpha_2$ balance different components in the combination of computational c	ost:
time complexity (forward and inverse computation time) and model size (parameter mem	ory
that mainly depends on NN depth in inverse problems).	
• $n_t = 0.01 n_t$ is used to control the speed of updating the trust-weight parameter $\beta$	bv
in sign and the second of aparating the trast weight parameter p	
evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found sta	ble
evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found statin adaptation.	ble
evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found statin adaptation.	ble
<ul> <li>evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found sta in adaptation.</li> <li>Specifically, we adopt trust-weight β to make the switching process stable and converge to an optimal statement of the statem</li></ul>	ble nal
<ul> <li>evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found sta in adaptation.</li> <li>Specifically, we adopt trust-weight β to make the switching process stable and converge to an optimiselection, which has been proven to improve online control in switching agents. β is updated based on the stability of the stab</li></ul>	ble nal sed
<ul> <li>evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found sta in adaptation.</li> <li>Specifically, we adopt trust-weight β to make the switching process stable and converge to an optime selection, which has been proven to improve online control in switching agents. β is updated ba on the difference in the combined evaluations of i-ResNet and DipDNN.</li> </ul>	ble nal sed
evaluating i-ResNet and DipDNN. In our implementation, 0.01 is empirically found statin adaptation. Specifically, we adopt trust-weight $\beta$ to make the switching process stable and converge to an optimate selection, which has been proven to improve online control in switching agents. $\beta$ is updated by the difference in the combined evaluations of i-ResNet and DipDNN.	ble nal sed

tational cost and acceptable accuracy, reflected by the trust-weight parameter favoring i-ResNet. For tasks requiring exact inversion, such as lossless data hiding or power system state estimation, the

algorithm shifts towards DipDNN due to its ability to provide a more accurate inverse, which results in a consistently lower evaluate function value for DipDNN.

A.7 NUMERICAL RESULTS

Table 1: Comparison of ResNet, I-ResNet, and DipDNN for MNIST Classification

Model	Config.	Acc.
DecNet	Conv	97.89%
Resivet	MLP	97.65%
I DogNot	Conv	96.32%
I-Residet	MLP	97.20%
DipDNN	MLP	97.69%

Table 2: Comparison of ResNet, I-ResNet, and DipDNN for CIFAR-10 Classification

Model	Block	Channels	Accuracy (CIFAR-10)
	12	16, 64, 256	89.64%
	12	64, 64, 64	89.21%
ResNet - Conv	12	16, 64, 256	90.79%
	39	12, 12, 12	88.06%
	39	12, 48, 192	90.69%
	12	64, 64, 64	83.22%
I-ResNet - Conv	12	16, 64, 256	87.79%
	39	12, 12, 12	84.69%
	39	16, 64, 256	90.00%
	39	12, 48, 192	90.26%
	12	12, 12, 12	84.29%
DipDNN - Conv	39	12, 12, 12	88.06%
_	39	12, 48, 192	90.69%
ResNet - MLP	12	-	66.97%
I-ResNet - MLP	12	-	66.69%

Table 3: Comparison of ResNet, i-ResNet, and DipDNN for physical system tasks.

(a) ResNet vs. DipDNN				
	ResNet	ResNet $k = 1$	ResNet $k = 2$	DipDNN
NS	0.00022	0.01498	0.0553	0.1258
<b>Power Flow</b>	5.79e-07	1.84e-05	3.19e-05	2.01e-06

	i-ResNet	i-ResNet $k = 1$	i-ResNet $k = 2$	DipDNN
NS	0.1635	7.7005	0.3343	0.1258
<b>Power Flow</b>	5.73e-05	2.08e-05	6.95e-05	2.01e-06

(b) i-ResNet vs. DipDNN





Figure 11: Forward and inverse performance of different models on Navier-Stokes flow dynamics.



Figure 12: Performance of convolutional models on Navier-Stokes flow dynamics.



Figure 13: Performance of different models on different image convert tasks.