

Smale's 18th Problem and the Barriers of Deep Learning

Matthew Colbrook

(University of Cambridge and École Normale Supérieure)

Smale's 18th problem*: *What are the limits of artificial intelligence?*

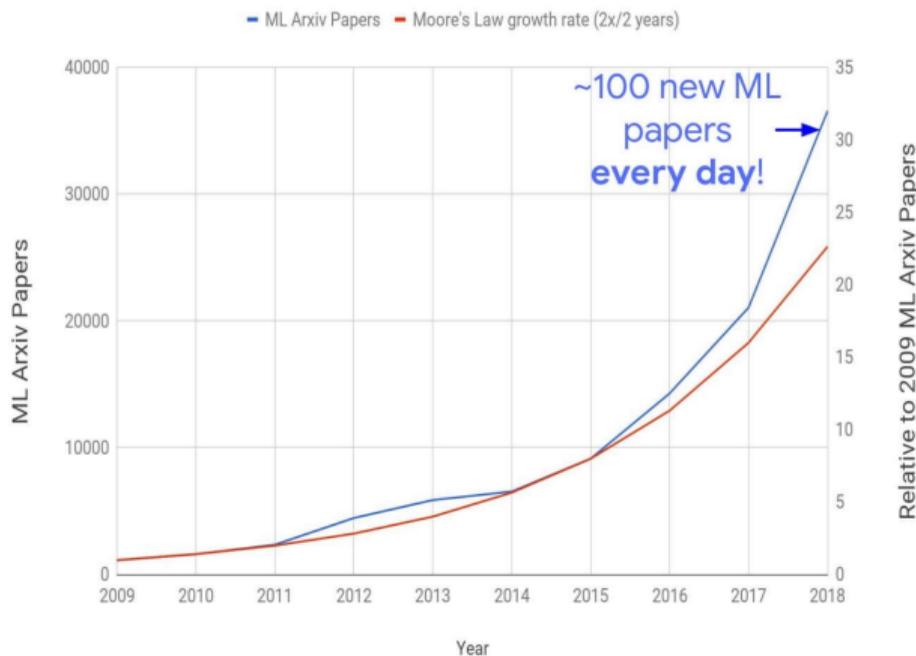
M. Colbrook, V. Antun, A. Hansen, “*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem*” (PNAS, 2022)

M. Colbrook, “*WARPd: A linearly convergent first-order method for inverse problems with approximate sharpness conditions*” (SIIMS, under revision)

*Steve Smale's list of problems for the 21st century (requested by Vladimir Arnold), inspired by Hilbert's list.

Interest in deep learning exponentially growing

Machine learning papers on arXiv



To keep up during first lockdown, would need to continually read a paper every 4 mins!

E.g., will AI replace standard algorithms in medical imaging?

nature > letters > article a nature research journal

MENU **nature** Search E-alert Submit Login

We'd like to understand how you use our websites in order to improve them. [Register your interest.](#)

Published: 22 March 2018

Image reconstruction by domain-transform manifold learning

Bo Zhu, Jeremiah Z. Liu, Stephen F. Cauley, Bruce R. Rosen & Matthew S. Rosen 

Nature 555, 487–492(2018) | [Cite this article](#)

17k Accesses | 235 Citations | 197 Altmetric | [Metrics](#)

Abstract

Image reconstruction is essential for imaging applications across the physical and life sciences, including optical and radar systems, magnetic resonance imaging, X-ray computed tomography, positron emission

You have full access to this article via University of Oslo Oslo University Hospital

[Download PDF](#) 

Editorial Summary

Machine learning improves image reconstruction

Reconstructing images from data, whether for medical or astronomical purposes, hinges on well-defined steps. The data sensor encodes an intermediate representation of the observed

[show all](#)

Claim: “superior immunity to noise and a reduction in reconstruction artefacts compared with conventional handcrafted reconstruction methods”.

Very strong confidence in deep learning

Forbes

Billionaires Innovation Leadership Money Consumer Industry

Turing Award And \$1 Million Given To 3 AI Pioneers

 **Nicole Martin** Contributor
AI & Big Data
I write about technology, data and privacy.

f
w
in



Winners of Turing Award sci.1016.1165

The Association for Computing Machinery (ACM) awarded Yoshua Bengio, Geoffrey Hinton and Yann LeCun with what many consider the "Nobel Prize of computing," for the innovations they've made in AI.

Cookies on Forbes

Geoffrey Hinton, The New Yorker, April 2017: "They should stop training radiologists now!"

Very strong confidence in deep learning

Forbes

Billionaires Innovation Leadership Money Consumer Industry

Turing Award And \$1 Million Given To 3 AI Pioneers

 **Nicole Martin** Contributor
AI & Big Data
I write about technology, data and privacy.

f
w
in



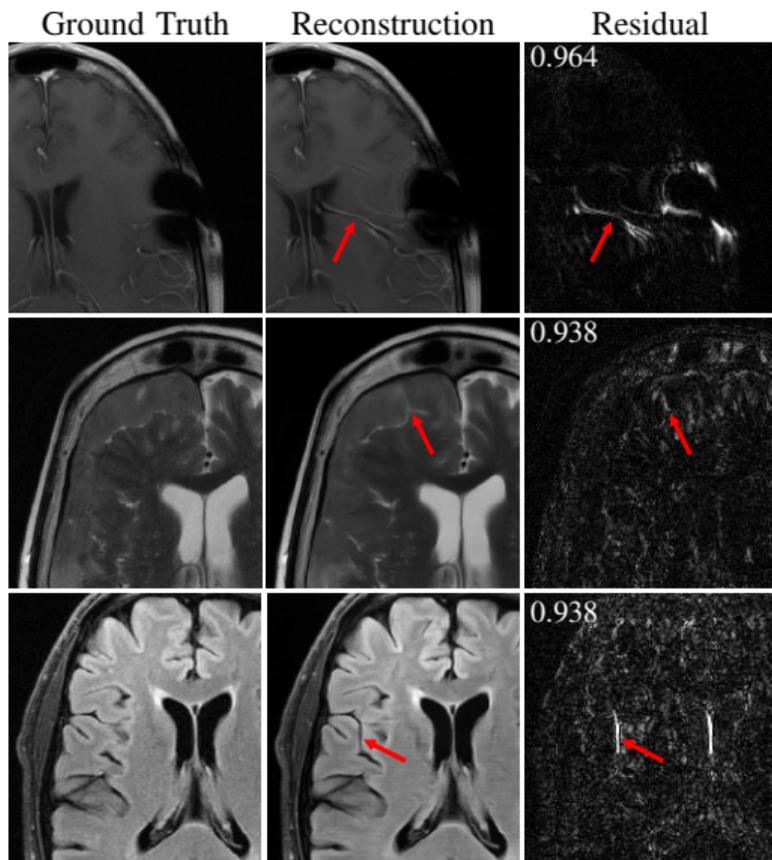
Winners of Turing Award scs 1016 1155

The Association for Computing Machinery (ACM) awarded Yoshua Bengio, Geoffrey Hinton and Yann LeCun with what many consider the "Nobel Prize of computing," for the innovations they've made in AI.

Cookies on Forbes

**Geoffrey Hinton, The New Yorker, April 2017: "They should stop training radiologists now!"
BUT ...**

AI hallucinations (Facebook and NYU's 2020 FastMRI challenge)

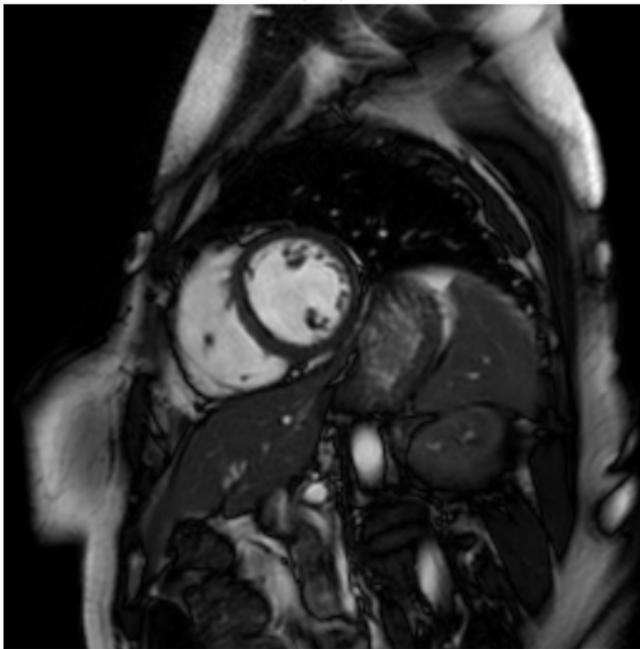


*"On AI, trust is a must, not a nice to have. High-risk AI systems will be subject to **strict obligations** before they can be put on the market: High level of **robustness, security and accuracy.**"*

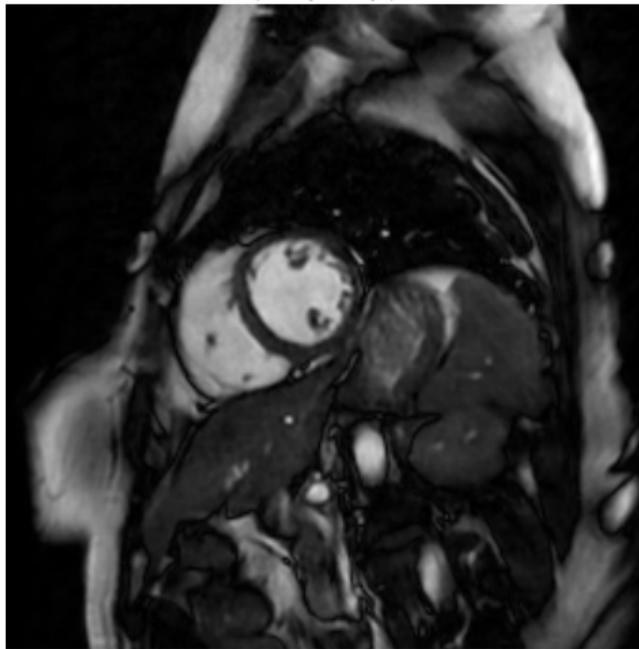
- Europ. Comm. outline for legal AI (April 2021).

Example of instabilities in inverse problems

$|x|$



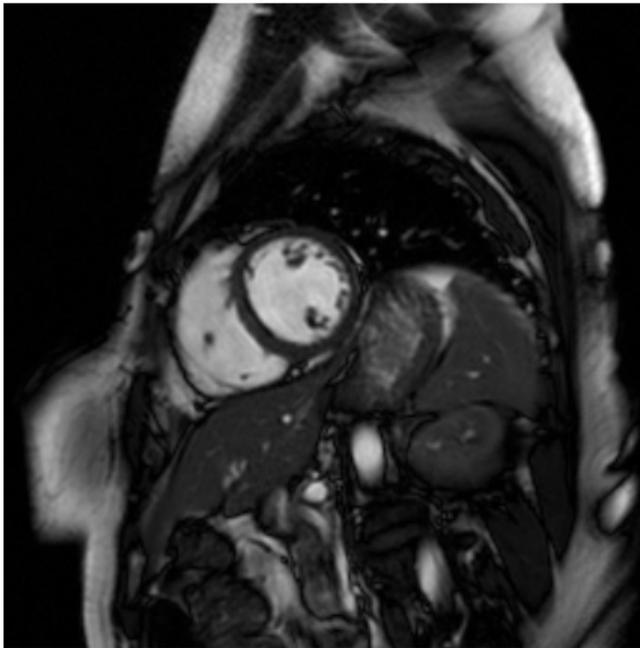
$|\Psi(Ax)|$



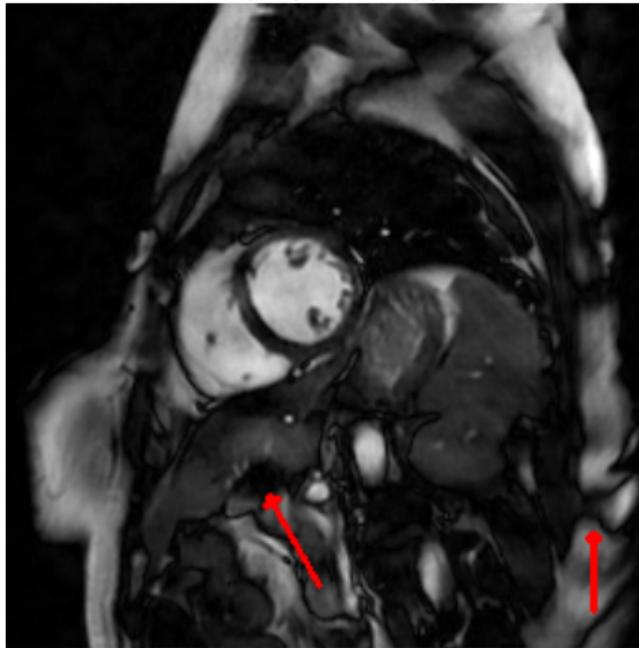
Network (33% subsampling) from: J. Schlemper, J. Caballero, J. V. Hajnal, A. Price and D. Rueckert, 'A deep cascade of convolutional neural networks for MR image reconstruction', in International conference on information processing in medical imaging, Springer, 2017, pp. 647–658.
Figures from: Antun, V., Renna, F., Poon, C., Adcock, B., & Hansen, A. C., 'On instabilities of deep learning in image reconstruction and the potential costs of AI'. Proc. Natl. Acad. Sci. USA, 2020.

Example of instabilities in inverse problems

$$|x + r_1|$$



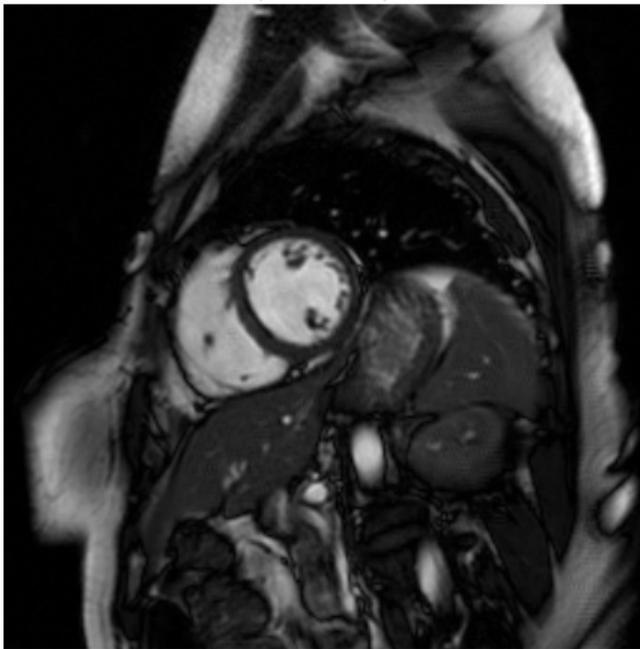
$$|\Psi(A(x + r_1))|$$



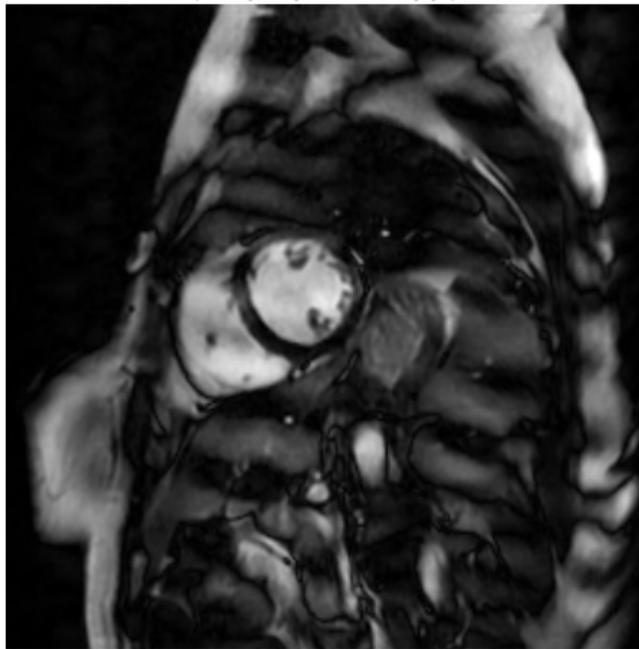
Network (33% subsampling) from: J. Schlemper, J. Caballero, J. V. Hajnal, A. Price and D. Rueckert, 'A deep cascade of convolutional neural networks for MR image reconstruction', in International conference on information processing in medical imaging, Springer, 2017, pp. 647–658.
Figures from: Antun, V., Renna, F., Poon, C., Adcock, B., & Hansen, A. C., 'On instabilities of deep learning in image reconstruction and the potential costs of AI'. Proc. Natl. Acad. Sci. USA, 2020.

Example of instabilities in inverse problems

$$|x + r_2|$$



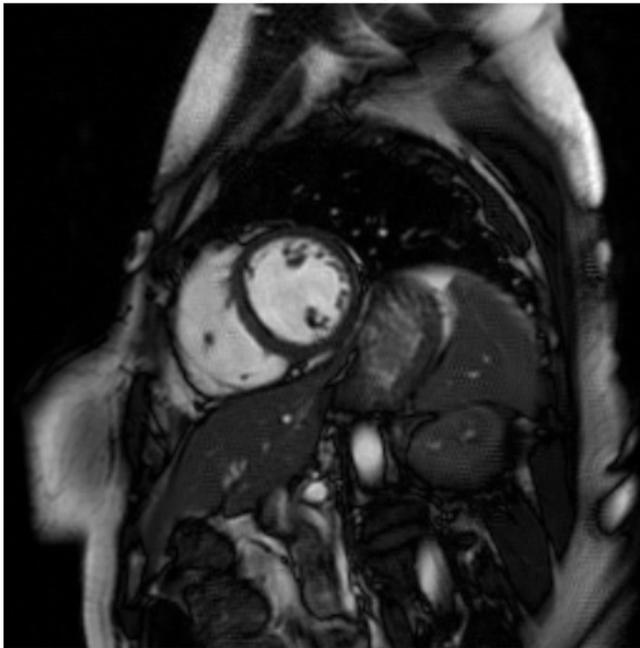
$$|\Psi(A(x + r_2))|$$



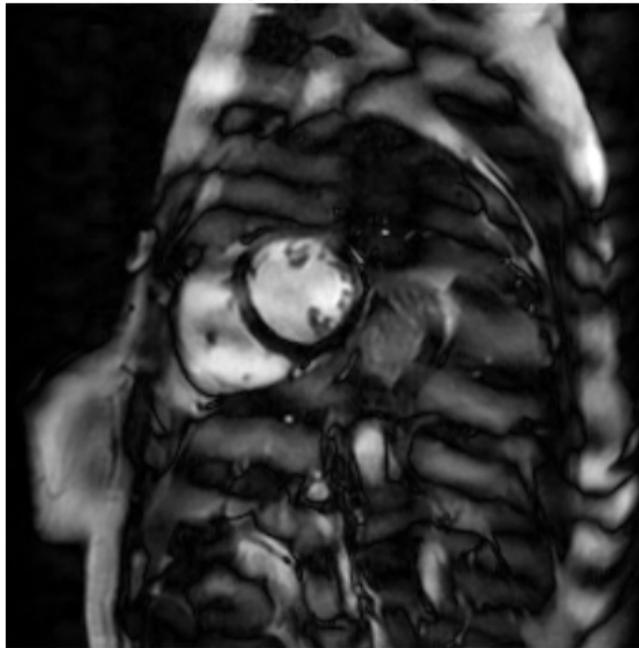
Network (33% subsampling) from: J. Schlemper, J. Caballero, J. V. Hajnal, A. Price and D. Rueckert, 'A deep cascade of convolutional neural networks for MR image reconstruction', in International conference on information processing in medical imaging, Springer, 2017, pp. 647–658.
Figures from: Antun, V., Renna, F., Poon, C., Adcock, B., & Hansen, A. C., 'On instabilities of deep learning in image reconstruction and the potential costs of AI'. Proc. Natl. Acad. Sci. USA, 2020.

Example of instabilities in inverse problems

$$|x + r_3|$$



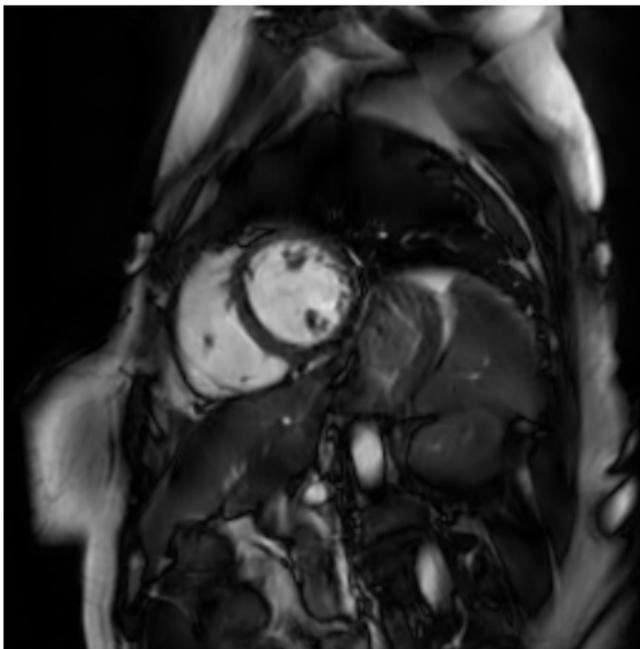
$$|\Psi(A(x + r_3))|$$



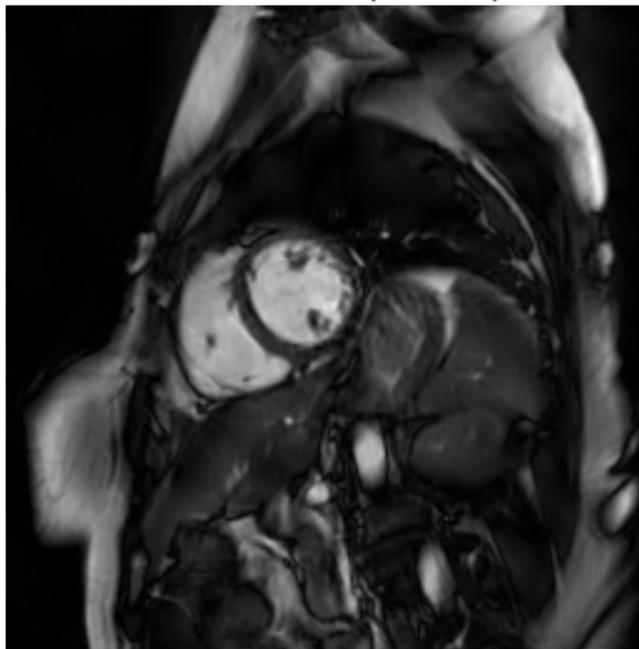
Network (33% subsampling) from: J. Schlemper, J. Caballero, J. V. Hajnal, A. Price and D. Rueckert, 'A deep cascade of convolutional neural networks for MR image reconstruction', in International conference on information processing in medical imaging, Springer, 2017, pp. 647–658.
Figures from: Antun, V., Renna, F., Poon, C., Adcock, B., & Hansen, A. C., 'On instabilities of deep learning in image reconstruction and the potential costs of AI'. Proc. Natl. Acad. Sci. USA, 2020.

Reconstruction using state-of-the-art standard methods

SoA from Ax



SoA from $A(x + r_3)$



Smale's 18th prob.: What are the limits of artificial intelligence?

*“Very often, the creation of a technological artifact precedes the science that goes with it. The steam engine was invented before thermodynamics. Thermodynamics was invented to explain the steam engine, essentially the **limitations** of it. What we are after is the equivalent of thermodynamics for intelligence.”*

— Yann LeCun (NYU, Facebook's chief AI scientist, Turing Award 2018)

*“2021 was the year in which the wonders of artificial intelligence stopped being a story. Many of this year's top articles grappled with the **limits of deep learning (today's dominant strand of AI).**”*

— IEEE Spectrum, 2021's Top Stories About AI (Dec. 2021)

Echoes of an old story

Hilbert's vision (start of 20th century): secure foundations for all mathematics.

- ▶ Mathematics written in a precise language.
- ▶ Completeness: all true math. statements can be proven.
- ▶ Consistency: no contradiction can be obtained.
- ▶ Decidability: algorithm for deciding truth of math. statements.



Hilbert's 10th problem: Provide an algorithm which, for any given polynomial equation with integer coefficients, can decide whether there is an integer-valued solution.

Foundations \Rightarrow better understanding, feasible directions for techniques, new methods, ...



Gödel (pioneer of **modern logic**) and Turing (pioneer of **modern computer science**):

- ▶ True statements in mathematics that cannot be proven!
- ▶ Computational problems that cannot be computed by an algorithm!

Hilbert's 10th problem: No such algorithm exists (1970, Matiyasevich).

A program for the foundations of DL and AI

A program determining the foundations/limitations of deep learning and AI is needed:

- ▶ Boundaries of methodologies.
- ▶ Universal/intrinsic boundaries (e.g., no algorithm can do it).

Key difference between existence and construction.

Two pillars of scientific computation:

- ▶ Stability
- ▶ Accuracy

GOAL of talk: Results in this direction for inverse problems.

Mathematical setup

Given $y = Ax + e$ recover $x \in \mathbb{C}^N$. $A \in \mathbb{C}^{m \times N}$, $m < N$ (e.g., MRI).

Outline:

- ▶ Paradox.
- ▶ Sufficient conditions and Fast Iterative REstarted NETworks (FIRENETs).
- ▶ Numerical examples (e.g., stability-accuracy trade-off).
- ▶ Approximate sharpness conditions and Weighted, Accelerated and Restarted Primal-dual (WARPd).

Can we train neural networks that solve (P_j) ?

$$\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} \quad \text{subject to} \quad \|Ax - y\|_{\ell^2} \leq \eta \quad (P_1)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}^2 \quad (P_2)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2} \quad (P_3)$$

Ξ = set of solutions.

Why P_j ?

- ▶ Avoid bizarre, unnatural & pathological mappings: (P_j) well-understood & well-used!
- ▶ Simpler solution map than inverse problem \Rightarrow stronger impossibility results.
- ▶ DL has also been used to speed up sparse regularization and tackle (P_j) .

The set-up

$$A \in \mathbb{C}^{m \times N} \text{ (modality)}, \quad \mathcal{S} = \{y_k\}_{k=1}^R \subset \mathbb{C}^m \text{ (samples)}, \quad R < \infty$$

In practice, A not known exactly or cannot be stored to infinite precision.

Assume access to: $\{y_{k,n}\}_{k=1}^R$ and A_n (rational approximations, e.g., floats) such that

$$\|y_{k,n} - y_k\| \leq 2^{-n}, \quad \|A_n - A\| \leq 2^{-n}, \quad \forall n \in \mathbb{N}.$$

Training set for $(A, \mathcal{S}) \in \Omega$:

$$\iota_{A, \mathcal{S}} := \{(y_{k,n}, A_n) \mid k = 1, \dots, R \text{ and } n \in \mathbb{N}\}.$$

In a nutshell: allow access to arbitrary precision training data.

The set-up

$$A \in \mathbb{C}^{m \times N} \text{ (modality)}, \quad \mathcal{S} = \{y_k\}_{k=1}^R \subset \mathbb{C}^m \text{ (samples)}, \quad R < \infty$$

In practice, A not known exactly or cannot be stored to infinite precision.

Assume access to: $\{y_{k,n}\}_{k=1}^R$ and A_n (rational approximations, e.g., floats) such that

$$\|y_{k,n} - y_k\| \leq 2^{-n}, \quad \|A_n - A\| \leq 2^{-n}, \quad \forall n \in \mathbb{N}.$$

Training set for $(A, \mathcal{S}) \in \Omega$:

$$\iota_{A, \mathcal{S}} := \{(y_{k,n}, A_n) \mid k = 1, \dots, R \text{ and } n \in \mathbb{N}\}.$$

In a nutshell: allow access to arbitrary precision training data.

Question: Given a collection Ω of (A, \mathcal{S}) , does there exist a neural network approximating Ξ (solution map of (P_j)), and can it be trained by an algorithm?

What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} \quad \text{subject to} \quad \|Ax - y\|_{\ell^2} \leq \eta \quad (P_1)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}^2 \quad (P_2)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2} \quad (P_3)$$

- (i) **Non-existence:** No neural network approximates Ξ .
- (ii)
- (iii)

What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} \quad \text{subject to} \quad \|Ax - y\|_{\ell^2} \leq \eta \quad (P_1)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}^2 \quad (P_2)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2} \quad (P_3)$$

- (i) ~~Non-existence:~~ No neural network approximates Ξ .
- (ii)
- (iii)

What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} \quad \text{subject to} \quad \|Ax - y\|_{\ell^2} \leq \eta \quad (P_1)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}^2 \quad (P_2)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2} \quad (P_3)$$

- (i) **Non-existence:** ~~No neural network approximates Ξ .~~
- (ii) **Non-trainable:** \exists a neural network that approximates Ξ , but it cannot be trained.
- (iii)

What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} \quad \text{subject to} \quad \|Ax - y\|_{\ell^2} \leq \eta \quad (P_1)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}^2 \quad (P_2)$$

$$\min_{x \in \mathbb{C}^N} \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2} \quad (P_3)$$

- (i) **Non-existence:** ~~No neural network approximates Ξ .~~
- (ii) **Non-trainable:** \exists a neural network that approximates Ξ , but it cannot be trained.
- (iii) **Not practical:** \exists a neural network that approximates Ξ , and an algorithm training it. However, the algorithm needs prohibitively many samples.

Paradox

Theorem

For (P_j) , $N \geq 2$ and $m < N$. Let $K \geq 3$ be a positive integer, $L \in \mathbb{N}$. Then there exists a **well-conditioned** class (condition numbers ≤ 1) Ω of elements (A, S) s.t. (Ω **fixed** in what follows):

Paradox

Theorem

For (P_j) , $N \geq 2$ and $m < N$. Let $K \geq 3$ be a positive integer, $L \in \mathbb{N}$. Then there exists a **well-conditioned** class (condition numbers ≤ 1) Ω of elements (A, S) s.t. (Ω **fixed** in what follows):

- (i) **There does not exist any algorithm** that, given a training set $\iota_{A,S}$, produces a neural network $\phi_{A,S}$ with

$$\min_{y \in \mathcal{S}} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-K}, \quad \forall (A, S) \in \Omega. \quad (1)$$

Furthermore, for any $p > 1/2$, **no probabilistic algorithm** can produce a neural network $\phi_{A,S}$ such that (1) holds with probability at least p .

Paradox

Theorem

For (P_j) , $N \geq 2$ and $m < N$. Let $K \geq 3$ be a positive integer, $L \in \mathbb{N}$. Then there exists a **well-conditioned** class (condition numbers ≤ 1) Ω of elements (A, S) s.t. (Ω fixed in what follows):

- (i) **There does not exist any algorithm** that, given a training set $\iota_{A,S}$, produces a neural network $\phi_{A,S}$ with

$$\min_{y \in S} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-K}, \quad \forall (A, S) \in \Omega. \quad (1)$$

Furthermore, for any $p > 1/2$, **no probabilistic algorithm** can produce a neural network $\phi_{A,S}$ such that (1) holds with probability at least p .

- (ii) **There exists an algorithm** that produces a neural network $\phi_{A,S}$ such that

$$\max_{y \in S} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-(K-1)}, \quad \forall (A, S) \in \Omega.$$

However, for any such algorithm (even probabilistic), $M \in \mathbb{N}$ and $p \in \left[0, 1 - \frac{1}{N+1-m}\right)$, there exists a training set $\iota_{A,S}$ such that for all $y \in S$,

$$\mathbb{P}\left(\inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} > 10^{-(K-1)} \text{ or size of training data needed} > M\right) > p.$$

Paradox

Theorem

For (P_j) , $N \geq 2$ and $m < N$. Let $K \geq 3$ be a positive integer, $L \in \mathbb{N}$. Then there exists a **well-conditioned** class (condition numbers ≤ 1) Ω of elements (A, S) s.t. (Ω fixed in what follows):

- (i) There **does not exist any algorithm** that, given a training set $\iota_{A,S}$, produces a neural network $\phi_{A,S}$ with

$$\min_{y \in S} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-K}, \quad \forall (A, S) \in \Omega. \quad (1)$$

Furthermore, for any $p > 1/2$, **no probabilistic algorithm** can produce a neural network $\phi_{A,S}$ such that (1) holds with probability at least p .

- (ii) There **exists an algorithm** that produces a neural network $\phi_{A,S}$ such that

$$\max_{y \in S} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-(K-1)}, \quad \forall (A, S) \in \Omega.$$

However, for any such algorithm (even probabilistic), $M \in \mathbb{N}$ and $p \in \left[0, 1 - \frac{1}{N+1-m}\right)$, there exists a training set $\iota_{A,S}$ such that for all $y \in S$,

$$\mathbb{P}\left(\inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} > 10^{-(K-1)} \text{ or size of training data needed} > M\right) > p.$$

- (iii) There **exists an algorithm** using only L training data from each $\iota_{A,S}$ that produces a neural network $\phi_{A,S}(y)$ such that

$$\max_{y \in S} \inf_{x^* \in \Xi(A,y)} \|\phi_{A,S}(y) - x^*\|_{\ell^2} \leq 10^{-(K-2)}, \quad \forall (A, S) \in \Omega.$$

In words ...

Nice classes Ω where stable and accurate neural networks exist. But:

- ▶ K digits: \nexists training algorithm for neural network.
- ▶ $K - 1$ digits: \exists training algorithm for neural network, but any such algorithm needs arbitrarily many training data.
- ▶ $K - 2$ digits: \exists training algorithm for neural network using L training samples.

Independent of neural network architecture - universal barrier.

Existence vs computation (universal approximation theorems **not** enough).

Conclusion: Theorems on existence of neural networks may have little to do with the neural networks produced in practice ...

Numerical example: fails with training methods

$\text{dist}(\Psi_{A_n}(y_n), \Xi(A, y))$	$\text{dist}(\Phi_{A_n}(y_n), \Xi(A, y))$	$\ A_n - A\ \leq 2^{-n}$ $\ y_n - y\ _{\ell^2} \leq 2^{-n}$	10^{-K}
0.2999690	0.2597827	$n = 10$	10^{-1}
0.3000000	0.2598050	$n = 20$	10^{-1}
0.3000000	0.2598052	$n = 30$	10^{-1}
0.0030000	0.0025980	$n = 10$	10^{-3}
0.0030000	0.0025980	$n = 20$	10^{-3}
0.0030000	0.0025980	$n = 30$	10^{-3}
0.0000030	0.0000015	$n = 10$	10^{-6}
0.0000030	0.0000015	$n = 20$	10^{-6}
0.0000030	0.0000015	$n = 30$	10^{-6}

Table: (Impossibility of computing the existing neural network to arbitrary accuracy).

Matrix $A \in \mathbb{C}^{19 \times 20}$ constructed from discrete cosine transform, $R = 8000$, solutions are 6-sparse. LISTA (learned iterative shrinkage thresholding algorithm) Ψ_{A_n} , and FIRENETs Φ_{A_n} . The table shows the shortest ℓ^2 distance between the output from the networks and the true minimizer of the problem $\min_{x \in \mathbb{C}^N} \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}$, for different values of n and K .

A paradox relevant to applications

IEEE Xplore FULL-TEXT LIBRARY IEEE STANDARDS IEEE Xplore

IEEE Spectrum FOR THE TECHNOLOGY ENGINEER

NEWS | ARTIFICIAL INTELLIGENCE

Some AI Systems May Be Impossible to Compute

New research suggests there are limitations to what deep neural networks can do

BY CHARLES Q. CHOI | 30 MAR 2022 | 4 MIN READ

For engineers



HOME NEWS RELEASES MULTIMEDIA MEETINGS

For scientists

NEWS RELEASE 17-MAR-2022

Mathematical paradoxes demonstrate the limits of AI

Peer-Reviewed Publication
UNIVERSITY OF CAMBRIDGE

Journal of the Society for Industrial and Applied Mathematics

siam news

Volume 55 Issue 4
May 2022

Privacy with Synthetic Data

For numerical analysts

Original Data
Differential Privacy
Vive Copula Model
Synthetic Data

Figure 1: Comparison of original and synthetic data. The synthetic data is generated using a generative model that has been trained on the original data. The synthetic data is indistinguishable from the original data.

Figure 2: Comparison of original and synthetic data. The synthetic data is generated using a generative model that has been trained on the original data. The synthetic data is indistinguishable from the original data.

Proving Existence Is Not Enough: Mathematical Paradoxes Unravel the Limits of Neural Networks in Artificial Intelligence

By Yogesh Anand, Matthew J. Coates, and Andre C. Courville

The impact of deep learning (DL) neural networks (NNs) and artificial intelligence (AI) over the last decade has been profound. Advances in computer vision, natural language processing, and other areas have led to significant improvements in accuracy and efficiency. However, the underlying mathematical principles of these models are often overlooked. In this paper, we explore the limits of NNs and AI by examining mathematical paradoxes that demonstrate the inherent limitations of these models. We show that NNs are fundamentally incapable of solving certain types of problems, and that AI is limited to a narrow range of tasks. These findings have important implications for the design and use of NNs and AI systems. We discuss the implications of these findings for the future of AI and NN research.

Figure 1: Synthetic Data

Figure 1 shows a comparison of original data (left) and synthetic data (right). The synthetic data is generated using a generative model that has been trained on the original data. The synthetic data is indistinguishable from the original data.

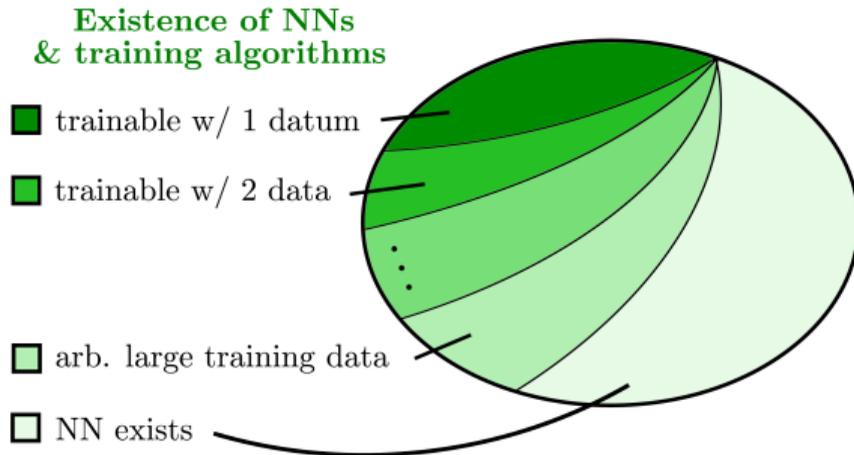
Figure 2: Synthetic Data

Figure 2 shows a comparison of original data (left) and synthetic data (right). The synthetic data is generated using a generative model that has been trained on the original data. The synthetic data is indistinguishable from the original data.

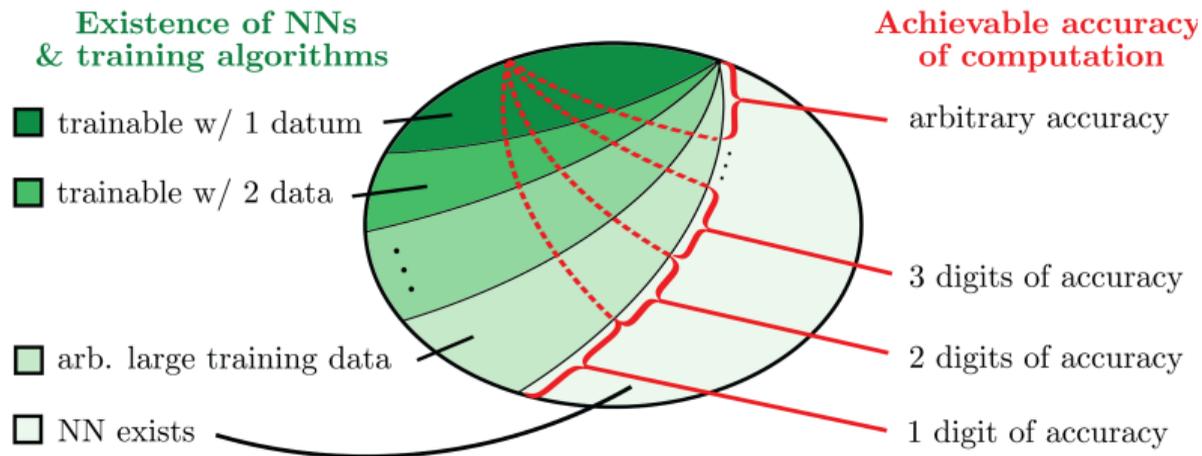
Figure 3: Synthetic Data

Figure 3 shows a comparison of original data (left) and synthetic data (right). The synthetic data is generated using a generative model that has been trained on the original data. The synthetic data is indistinguishable from the original data.

The world of neural networks



The world of neural networks



Need: Classification theory saying what can/cannot be done.

Example:

$$\hat{x} \in \operatorname{argmin} f(x), \quad f^* = \min f(x)$$

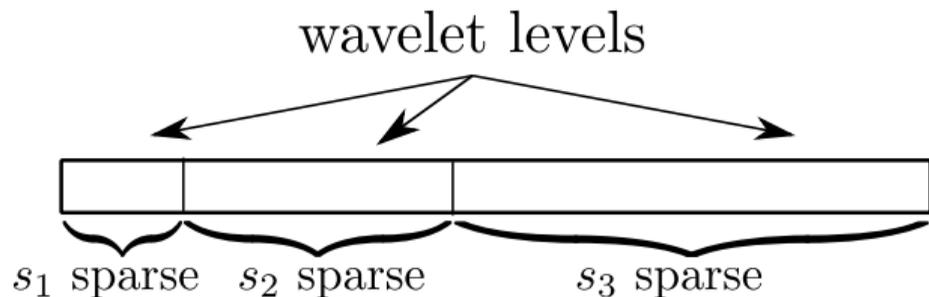
Problem: $f(x) \leq f^* + \epsilon$ does not in general imply x is close to set of minimizers.

Question: Can we find 'good' input classes where

$$f(x) \leq f^* + \epsilon \implies \inf_{\hat{x} \in \operatorname{argmin} f(x)} \|x - \hat{x}\| \lesssim \epsilon?$$

We shall see that the answer is yes!

State-of-the-art model for sparse regularisation



$\mathbf{M} = (M_1, \dots, M_r) \in \mathbb{N}^r$ and $\mathbf{s} = (s_1, \dots, s_r) \in \mathbb{Z}_{\geq 0}^r$. $x \in \mathbb{C}^N$ is (\mathbf{s}, \mathbf{M}) -sparse in levels if

$$|\text{supp}(x) \cap \{M_{k-1} + 1, \dots, M_k\}| \leq s_k, \quad k = 1, \dots, r.$$

Denote set of (\mathbf{s}, \mathbf{M}) -sparse vectors by $\Sigma_{\mathbf{s}, \mathbf{M}}$, define

$$\sigma_{\mathbf{s}, \mathbf{M}}(x)_{\ell^1} = \inf \{ \|x - z\|_{\ell^1} : z \in \Sigma_{\mathbf{s}, \mathbf{M}} \}.$$

The robust nullspace property

Definition: $A \in \mathbb{C}^{m \times N}$ satisfies the **robust null space property in levels (rNSPL)** of order (\mathbf{s}, \mathbf{M}) with constants $\rho \in (0, 1)$ and $\gamma > 0$ if for any (\mathbf{s}, \mathbf{M}) support set Δ ,

$$\|x_{\Delta}\|_{\ell^2} \leq \frac{\rho \|x_{\Delta^c}\|_{\ell^1}}{\sqrt{r(s_1 + \dots + s_r)}} + \gamma \|Ax\|_{\ell^2}, \quad \forall x \in \mathbb{C}^N.$$

Objective function: $f(x) = \lambda \|x\|_{\ell^1} + \|Ax - y\|_{\ell^2}$

$$\begin{aligned} \text{rNSPL} \Rightarrow \|z - x\|_{\ell^2} &\lesssim \underbrace{\sigma_{\mathbf{s}, \mathbf{M}}(x)_{\ell^1} + \|Ax - y\|_{\ell^2}}_{\text{"small"}} \\ &\quad + \underbrace{(\lambda \|z\|_{\ell^1} + \|Az - y\|_{\ell^2} - \lambda \|x\|_{\ell^1} - \|Ax - y\|_{\ell^2})}_{f(z) - f(x) \text{ objective function difference}}, \end{aligned}$$

In a nutshell: control $\|z - x\|_{\ell^2}$ by $f(z) - f(x)$, up to small approximation term.

Fast Iterative REstarted NETworks (FIRENETs)

Simplified version of Theorem: *We provide an algorithm such that:*

Input: *Sparsity parameters (\mathbf{s}, \mathbf{M}) , $A \in \mathbb{C}^{m \times N}$ satisfying the rNSPL with constants $0 < \rho < 1$ and $\gamma > 0$, $n \in \mathbb{N}$ and positive $\{\delta, b_1, b_2\}$.*

Output: *A neural network ϕ_n with $\mathcal{O}(n)$ layers and width $2(N + m)$ such that:*

For any $x \in \mathbb{C}^N$ and $y \in \mathbb{C}^m$ with

$$\underbrace{\sigma_{\mathbf{s}, \mathbf{M}}(x)_{\ell^1}}_{\text{distance to sparse in levels vectors}} + \underbrace{\|Ax - y\|_{\ell^2}}_{\text{noise of measurements}} \lesssim \delta, \quad \|x\|_{\ell^2} \lesssim b_1, \quad \|y\|_{\ell^2} \lesssim b_2,$$

*we have the following **stable** and **exponential convergence** guarantee in n*

$$\|\phi_n(y) - x\|_{\ell^2} \lesssim \delta + e^{-n}.$$

Demonstration of convergence

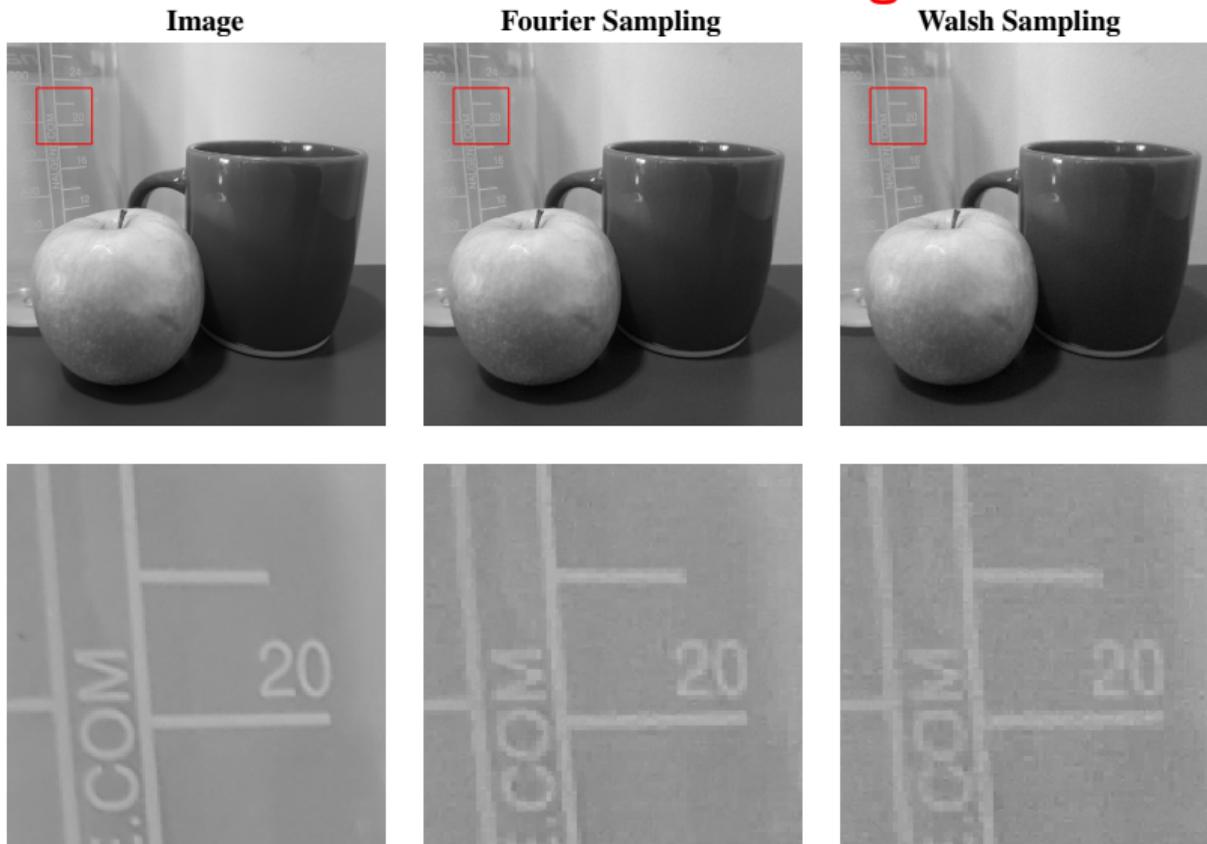
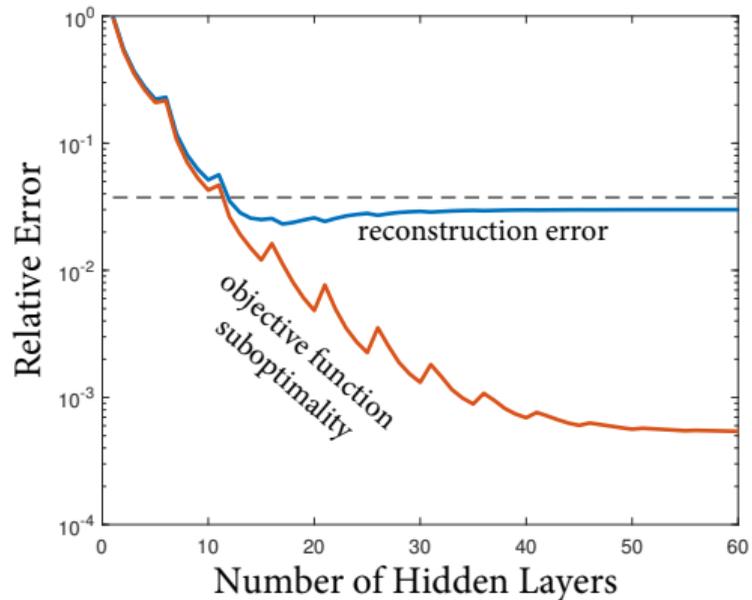


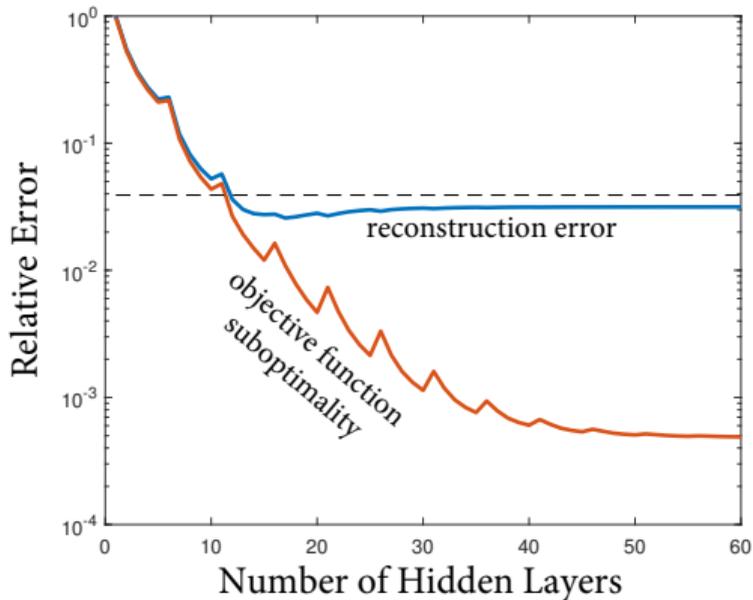
Figure: Images corrupted with 2% Gaussian noise and reconstructed using 15% sampling.

Demonstration of convergence

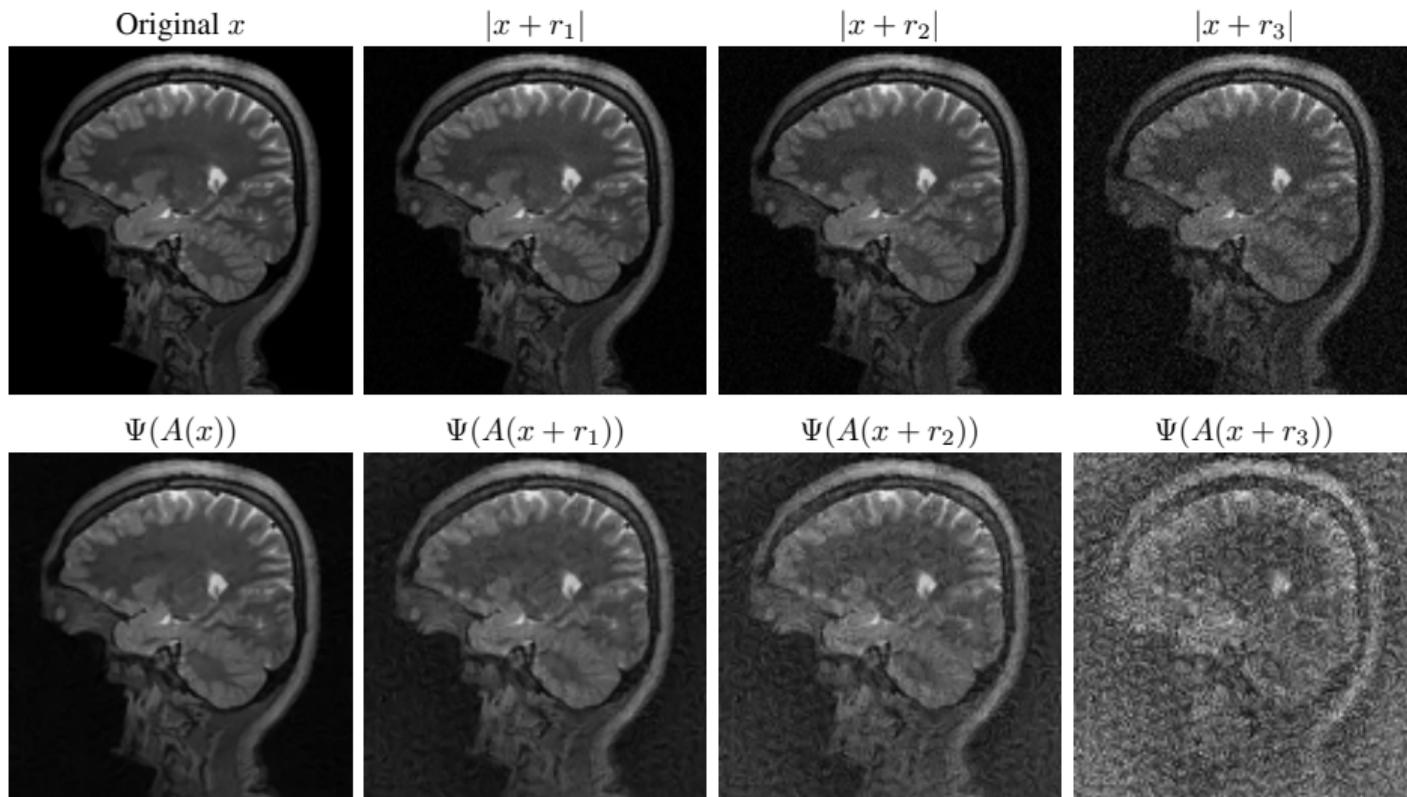
Convergence, Fourier Sampling



Convergence, Walsh Sampling

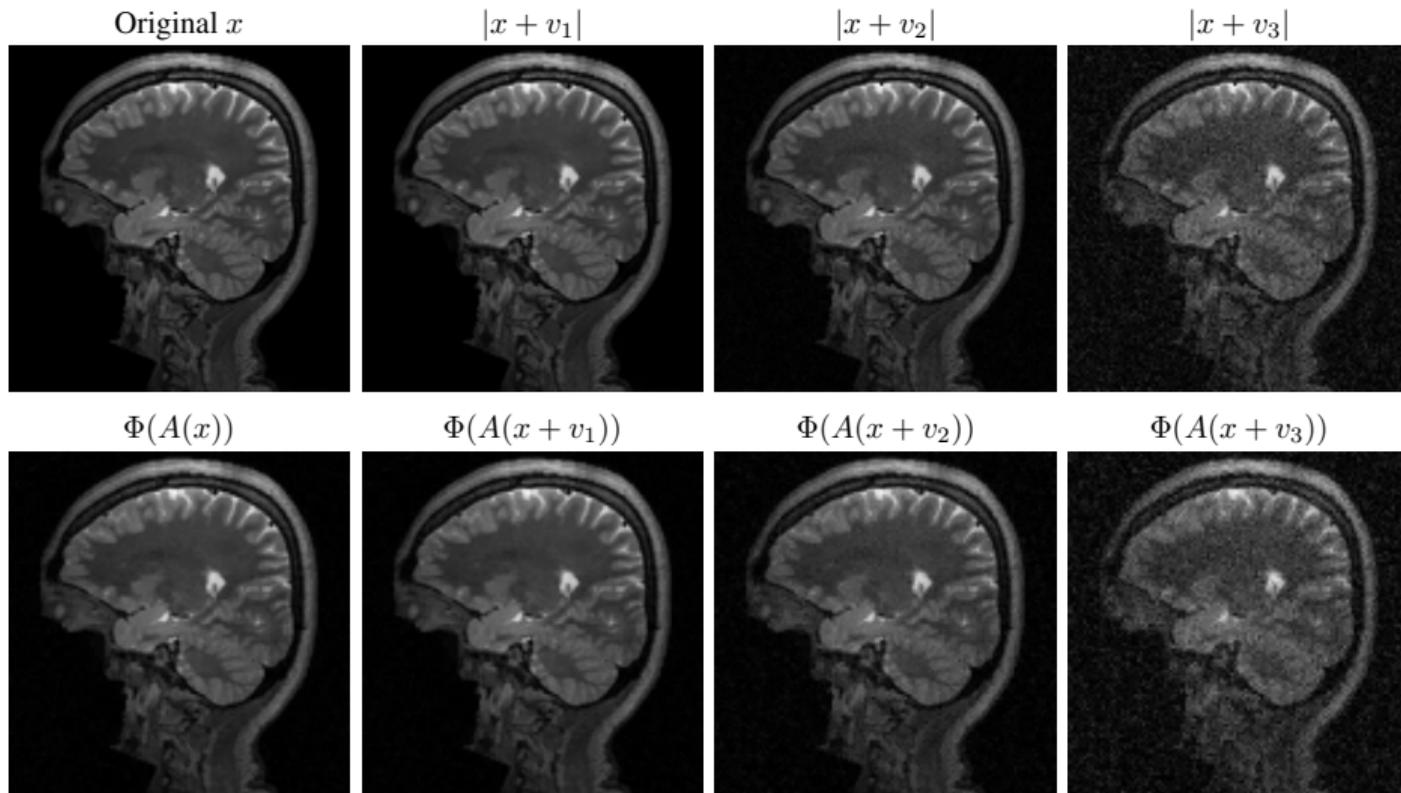


Stable? AUTOMAP \times



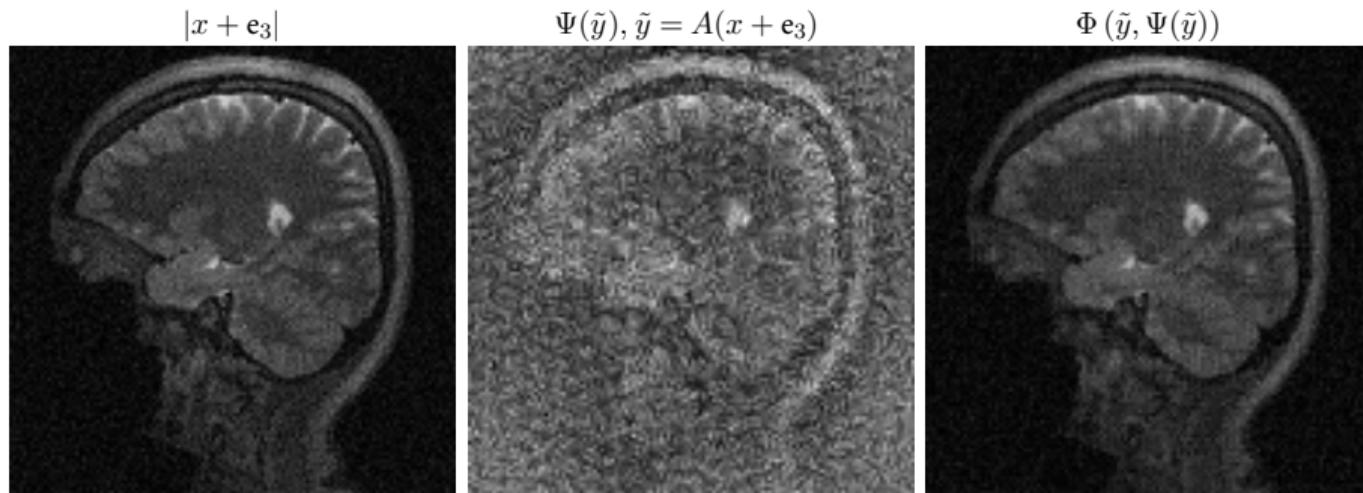
- V. Antun et al. "On instabilities of deep learning in image reconstruction and the potential costs of AI," PNAS, 2021.
- B. Zhu et al. "Image reconstruction by domain-transform manifold learning," Nature, 2018.

Stable? FIRENETs ✓



-
- M. Colbrook, V. Antun, A. Hansen, "The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem," PNAS, 2022.

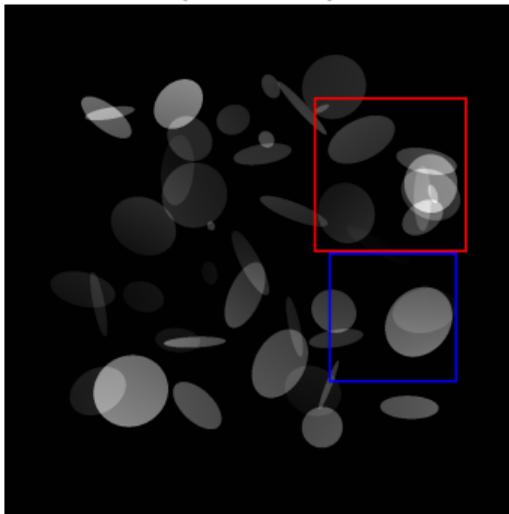
Adding FIRENET layers stabilizes AUTOMAP



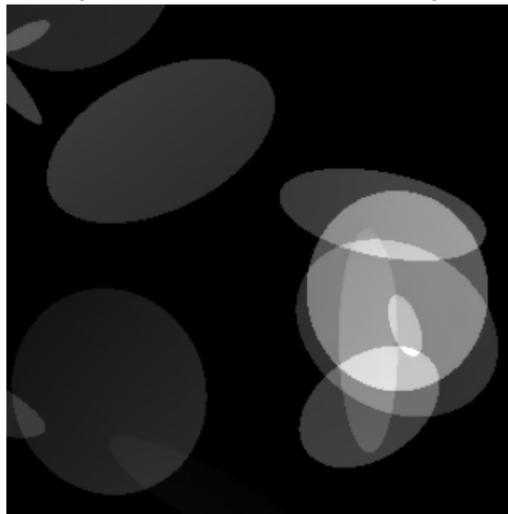
-
- M. Colbrook, V. Antun, A. Hansen, "The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem," PNAS, 2022.

Stability vs. accuracy tradeoff

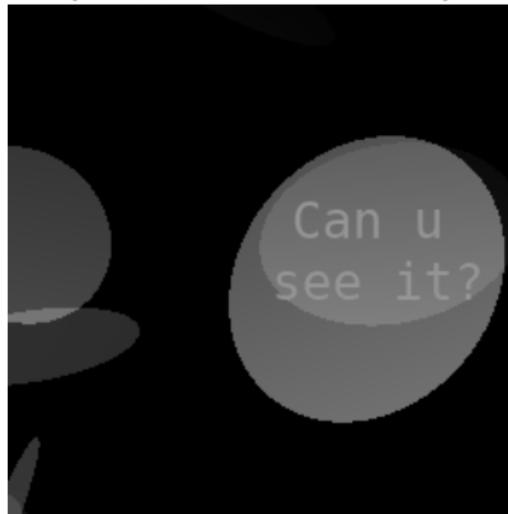
Original x
(full size)



Original
(cropped, red frame)



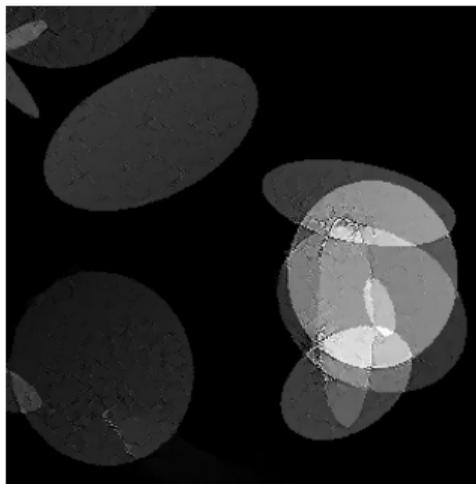
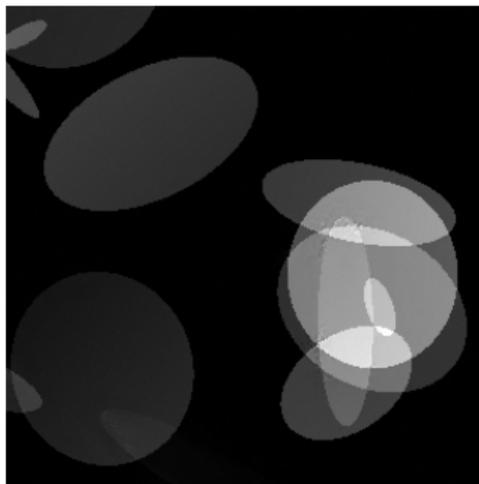
Original + detail ($x + h_1$)
(cropped, blue frame)



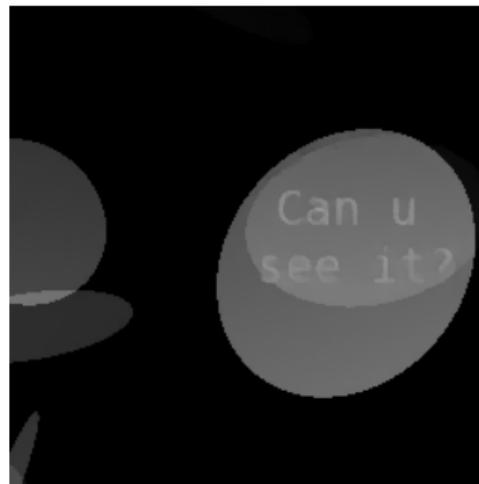
-
- M. Colbrook, V. Antun, A. Hansen, "The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem," PNAS, 2022.

U-net trained without noise

Orig. + worst-case noise Rec. from worst-case noise

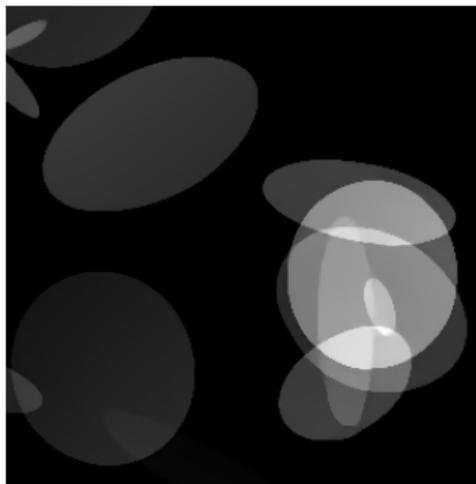
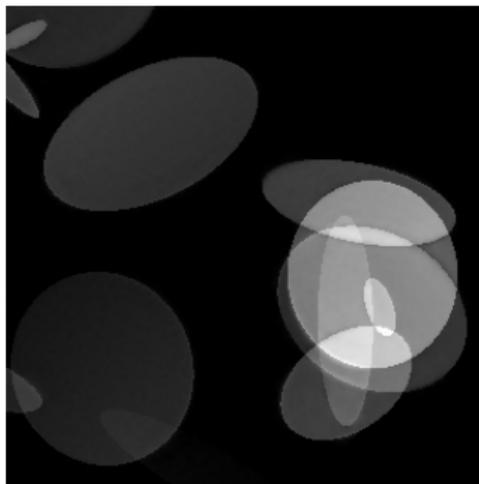


Rec. of detail

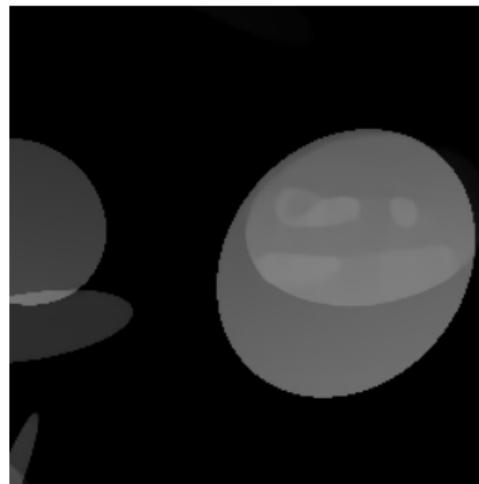


U-net trained with noise

Orig. + worst-case noise Rec. from worst-case noise

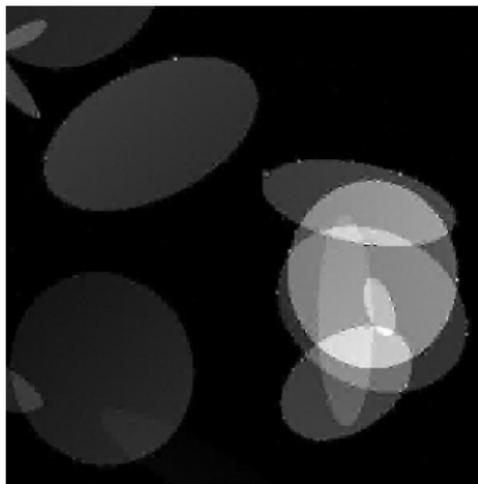
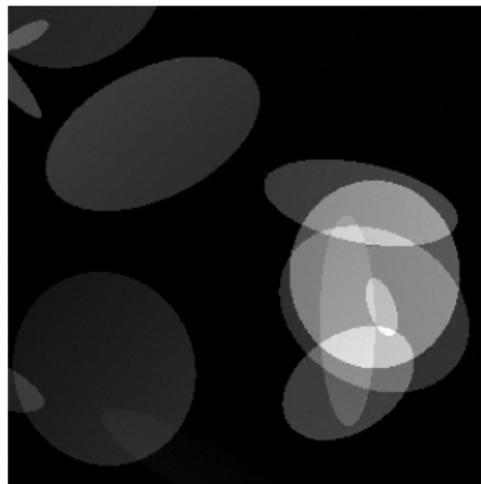


Rec. of detail



FIRENET

Orig. + worst-case noise Rec. from worst-case noise



Rec. of detail



Broader framework: approximate sharpness conditions

Problem: Given $y = Ax + e \in \mathbb{C}^m$, recover $x \in \mathbb{C}^N$.

Optimization: $\min_{x \in \mathbb{C}^N} \mathcal{J}(x) + \|Bx\|_{\ell^1}$ s.t. $\|Ax - y\|_{\ell^2} \leq \epsilon$, $B \in \mathbb{C}^{q \times N}$.

Assume: $\|\hat{x} - x\|_{\ell^2} \leq C_1 \left[\underbrace{\mathcal{J}(\hat{x}) + \|B\hat{x}\|_{\ell^1} - \mathcal{J}(x) - \|Bx\|_{\ell^1}}_{\text{objective function difference}} + C_2 \left(\underbrace{\|A\hat{x} - y\|_{\ell^2} - \epsilon}_{\text{feasibility gap}} + \underbrace{c(x, y)}_{\text{approx. term}} \right) \right]$.

Examples: Sparse vector recovery, low-rank matrix recovery, matrix completion, ℓ^1 -analysis problems, TV minimization, mixed regularization problems, ...

Simplified version of Theorem: Let $\delta > 0$. We provide a neural network ϕ of depth $\mathcal{O}(\log(\delta^{-1}))$ and width $\mathcal{O}(N + m + q)$ such that for all $(x, y) \in \mathbb{C}^N \times \mathbb{C}^m$

$$\|Ax - y\|_{\ell^2} \leq \epsilon \text{ and } c(x, y) \leq \delta \quad \Rightarrow \quad \|\phi(y) - x\|_{\ell^2} \lesssim \delta.$$

Weighted, Accelerated and Restarted Primal-dual (WARPd)

- ▶ Primal-dual iterations starting at x_0 :

$$\underbrace{\mathcal{J}(X_k) + \|BX_k\|_{\ell^1} - \mathcal{J}(x) - \|Bx\|_{\ell^1} + C_2 (\|AX_k - b\|_{\ell^2} - \epsilon)}_{=: G(X_k)} \leq \frac{1}{k} \left(\frac{\|x_0 - x\|_{\ell^2}^2}{\tau_1} + \frac{C_2^2 + q}{\tau_2} \right). \quad (2)$$

- ▶ Assumption implies $\|X_k - x\|_{\ell^2} \leq C_1(G(X_k) + \delta)$, controls RHS of (2) upon restart.
- ▶ Reweight and optimize parameters for map H_k using k iterations s.t.

$$G(x_0) \leq \alpha_0 \Rightarrow G(H_k(x_0)) \leq C(\delta + \alpha_0)/k$$

- ▶ Restart when $C/k \leq e^{-1}$ (optimal). \tilde{X}_p after p restarts:

$$G(\tilde{X}_p) \leq e^{-1}(\delta + e^{-1}(\delta + \dots + e^{-1}(\delta + \alpha_0))) = (e^{-1} + e^{-2} + \dots + e^{-p})\delta + e^{-p}\alpha_0 \lesssim \delta + e^{-p}.$$

- ▶ Apply the assumption to get $\|\tilde{X}_p - x\|_{\ell^2} \lesssim \delta + e^{-p}$.

Remarks:

- ▶ Unrolled as neural networks. **NB:** Naive unrolling gives slow $\mathcal{O}(\delta + p^{-1})$ convergence.
- ▶ Stability w.r.t. input (inherent) and execution (numerical).

- M. Colbrook "WARPd: A linearly convergent first-order method for inverse problems with approximate sharpness conditions."
- A. Chambolle, T. Pock, "A first-order primal-dual algorithm for convex problems with applications to imaging," J Math Imaging Vis, 2011.
- V. Monga, Y. Li, Y. C. Eldar, "Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing," IEEE Signal Process. Mag., 2021.

A final example with different regularizers

A is a DFT, 15% subsampled according to an inverse square law (optimal for TV). Measurements are corrupted with 5% Gaussian noise.

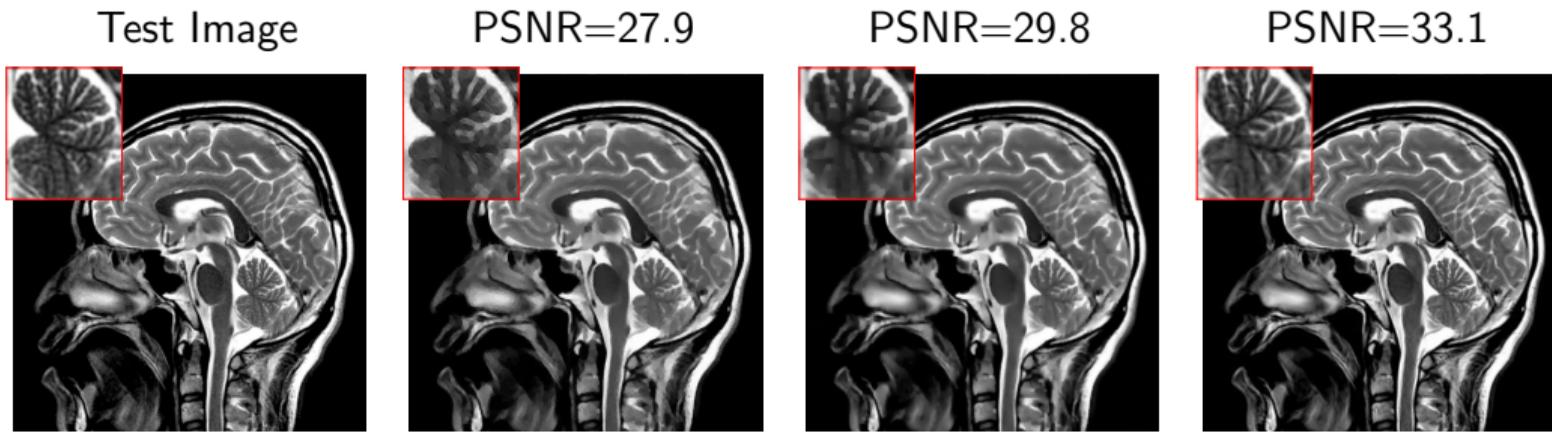


Figure: Middle-left: Converged reconstruction using TV. Middle-right: Converged reconstruction using TGV. Right: Reconstruction using (adaptively adjusted weighted) shearlets and TGV, after 25 iterations. All reconstructions were computed using WARPd.

WARPd can easily handle complicated mixed regularization problems.

$$\min_{x \in \mathbb{C}^N} \|WD^*x\|_{\ell^1} + \text{TGV}_\alpha^2(x) \quad \text{s.t.} \quad \|Ax - b\|_{\ell^2} \leq \epsilon,$$

Concluding remarks

There is a **need for foundations** in AI/deep learning.

- ▶ ‘Nice’ inverse problems where stable & accurate neural network exists but cannot be trained!
- ▶ Existence of training algorithm depends on desired accuracy. $\forall K \in \mathbb{Z}_{\geq 3}, \exists$ classes s.t.:
 - (i) Algorithms may compute neural networks to $K - 1$ digits of accuracy, but not K .
 - (ii) Achieving $K - 1$ digits of accuracy requires arbitrarily many training data.
 - (iii) Achieving $K - 2$ correct digits requires only one training datum.
- ▶ Under *specific conditions*, algorithms can train stable and accurate neural networks. E.g., **FIRENETs** achieve exponential convergence & withstand adversarial attacks.
- ▶ Trade-off between stability and accuracy in deep learning.

- ▶ **WARPd** provides accelerated recovery under an approximate sharpness condition.
- ▶ Quantities controlling recovery also provide explicit approximate sharpness constants.
- ▶ $\text{WARPd}_{\text{unrolled}} \Rightarrow \text{FIRENETs}$.

Question: How do we optimally traverse the stability & accuracy trade-off?