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

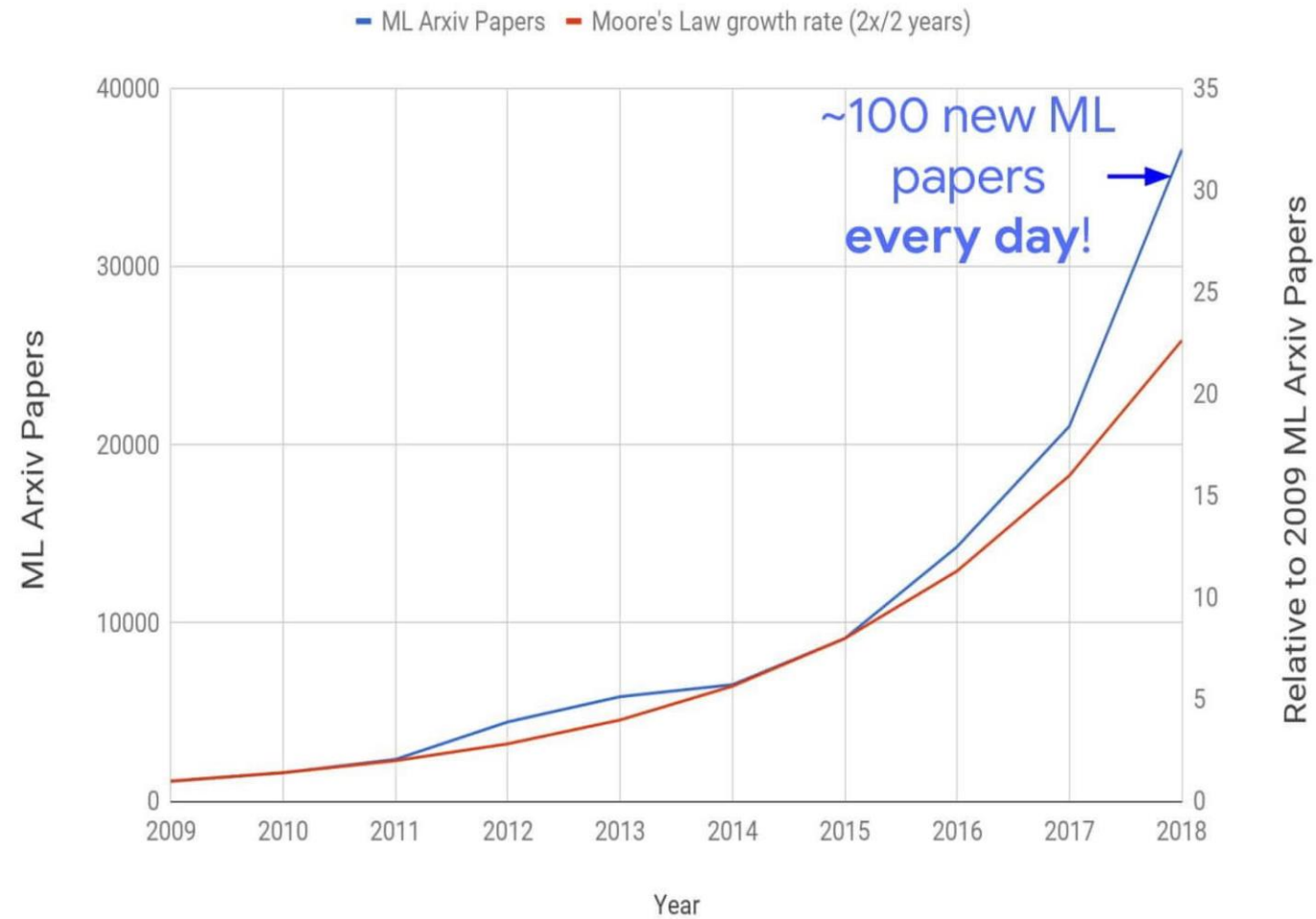**Matthew Colbrook**

University of Cambridge

> **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)

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

http://www.damtp.cam.ac.uk/user/mjc249/home.html: slides, papers, and code
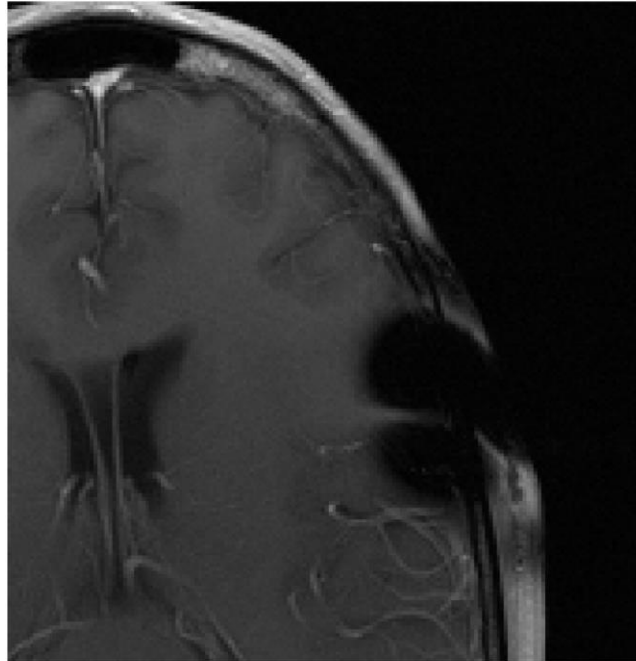
# A fun stat!



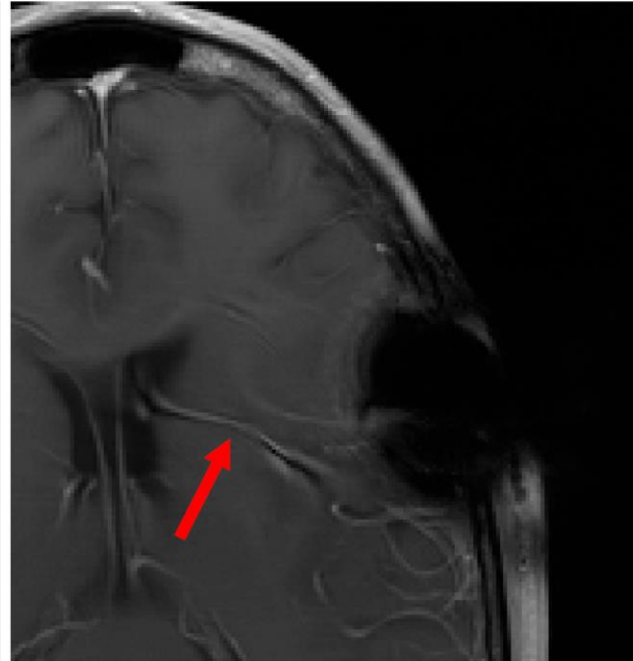To keep up during first lockdown, would need to continually read a paper every 4 mins!

# Problem: hallucinations and instability



**Hallucinations in image reconstruction**
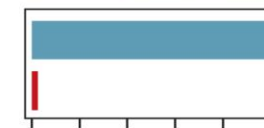Original image — AI reconstruction

**Instabilities in medical diagnosis**
Original Mole — Perturbed Mole

Benign / Malignant — Model confidence
Benign / Malignant — Model confidence

"**AI generated hallucination**", from Facebook and NYU's *FastMRI challenge* 2020

From Finlayson et al., "*Adversarial attacks on medical machine learning*," **Science**, 2019.

# When can we make AI robust and trustworthy?

# Smale's 18th problem: "What are the limits of AI?"

*"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

*"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)

# Example of the limits of deep learning

**Paradox:** "Nice" linear inverse problems where a *stable* and *accurate* neural network for image reconstruction <u>exists</u>, but it <u>can never be trained</u>!

E.g., suppose we want to solve (holds for much more general problems)

$$\min_{x \in \mathbb{C}^N} \|x\|_{l^1} + \lambda \|Ax - y\|_{l_2}^2$$

$$A \in \mathbb{C}^{m \times N} \ (\text{modality}, m < N), \qquad S = \{y_K\}_{K=1}^{R} (\text{samples})$$

Arises when given $y \approx Ax + e$.

Enforce condition numbers bounded by 1.

# Data

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

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

Assume access to $\{y_{n,k}\}_{k=1}^{R}$ and $A_n$ (rational approx, e.g., floats) such that
$$\|y_{n,k} - y_k\| \leq 2^{-n}, \qquad \|A_n - A\| \leq 2^{-n}, \qquad n \in \mathbb{N}.$$

Training set for $(A, S) \in \Omega$:
$$\iota_{A,S} = \{(y_{n,k}, A_n) : k = 1, \ldots, R \text{ and } n \in \mathbb{N}\}.$$

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

**Question:** Given a collection $\Omega$ of $(A, S)$, does there exist a neural network approximating the solution map, and can it be trained by an algorithm?

# What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{l^1} + \lambda \|Ax - y\|_{l_2}^2$$

What could go wrong?

1. **Non-existence:** No neural network approximates solution map.

# What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{l^1} + \lambda \|Ax - y\|_{l_2}^2$$

What could go wrong?

1. ~~**Non-existence:** No neural network approximates solution map.~~

# What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{l^1} + \lambda \|Ax - y\|_{l_2}^2$$

What could go wrong?

1. ~~**Non-existence:** No neural network approximates solution map.~~

2. **Non-trainable:** ∃ a neural network that approximates solution map, but it cannot be trained.

# What could go wrong?

$$\min_{x \in \mathbb{C}^N} \|x\|_{l^1} + \lambda \|Ax - y\|_{l_2}^2$$

What could go wrong?

1. ~~**Non-existence:** No neural network approximates solution map.~~

2. **Non-trainable:** ∃ a neural network that approximates solution map, but it cannot be trained.

3. **Not practical:** ∃ a neural network that approximates solution map, and an algorithm training it. However, the algorithm needs prohibitively many samples.

# Example of the limits of deep learning

**Paradox:** "Nice" linear inverse problems where a *stable* and *accurate* neural network for image reconstruction <u>exists</u>, but it <u>can never be trained</u>!

<span style="background:yellow">**Theorem**</span>: Pick positive integers $n \geq 3$ and $M$. Class of problems such that:
- (**Not trainable**) No algorithm (even random) can train a neural network with $\boldsymbol{n}$ **digits** of accuracy over the dataset with probability greater than 1/2.
- (**Not practical**) $\boldsymbol{n-1}$ **digits** of accuracy possible over the dataset, but any training algorithm requires **arbitrarily large training data**.
- (**Trainable and practical**) $\boldsymbol{n-2}$ digits of accuracy possible over the dataset via training algorithm using $\boldsymbol{M}$ **training data**.

Holds for any architecture, any precision of training data.

$\Longrightarrow$ **Classification theory telling us what can and cannot be done**

- C., Antun, Hansen, "*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem,*" **PNAS**, 2022.
- Antun, C., Hansen, *"Proving Existence Is Not Enough: : Mathematical Paradoxes Unravel the Limits of Neural Networks in Artificial Intelligence,"* **SIAM News,** May 2022.
- Choi, *"Some AI Systems May Be Impossible to Compute,"* **IEEE Spectrum,** March 2022.

# 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 | $n = 30$ | $10^{-6}$ |

$A \in \mathbb{C}^{19 \times 20}$ from discrete cosine transform, $R = 8000$, solutions 6-sparse. LISTA (learned iterative shrinkage thresholding algorithm) $\Psi_{A_n}$ and FIRENETs $\Phi_{A_n}$. The table shows the shortest $l_2$ distance between the output and the true minimizer of the problem for different values of $n, K$.

# A paradox relevant to applications

# The world of neural networks



Existence of NNs
& training algorithms

- ■ trainable w/ 1 datum
- ■ trainable w/ 2 data
- ■ arb. large training data
- □ NN exists

**Given a problem and conditions, where does it sit in this diagram?**

# The world of neural networks



**Given a problem and conditions, where does it sit in this diagram?**

# Example counterpart theorem

**Certain conditions:** <u>stable</u> neural networks <u>trained</u> with <u>exponential accuracy</u>. E.g., *approximate Łojasiewicz-type inequality*:

$$(1) \quad \min_{x \in \mathbb{C}^N} f(x) \quad \text{s.t.} \quad \|Ax - y\| \leq \varepsilon$$

$$\text{dist}(x, \text{solution}) \leq \alpha([f(x) - f^*] + [\|Ax - y\| - \varepsilon] + \delta)$$

**F**ast **I**terative **RE**started **NET**works (FIRENETs)
(unrolled primal-dual with novel restart scheme)

**Theorem**: Training algorithm that, under above assumption, produces *stable* neural networks $\varphi_n$ of width $O(N)$, depth $O(n)$, guaranteed worst bound

$$\text{dist}(\varphi_n(y), \text{solution}) \lesssim e^{-n} + \delta$$

- C., Antun, Hansen, "*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem*," **PNAS**, 2022.
- C., "*WARPd: A linearly convergent first-order method for inverse problems with approximate sharpness conditions*," **SIIMS**, 2022.

# Demonstration of convergence

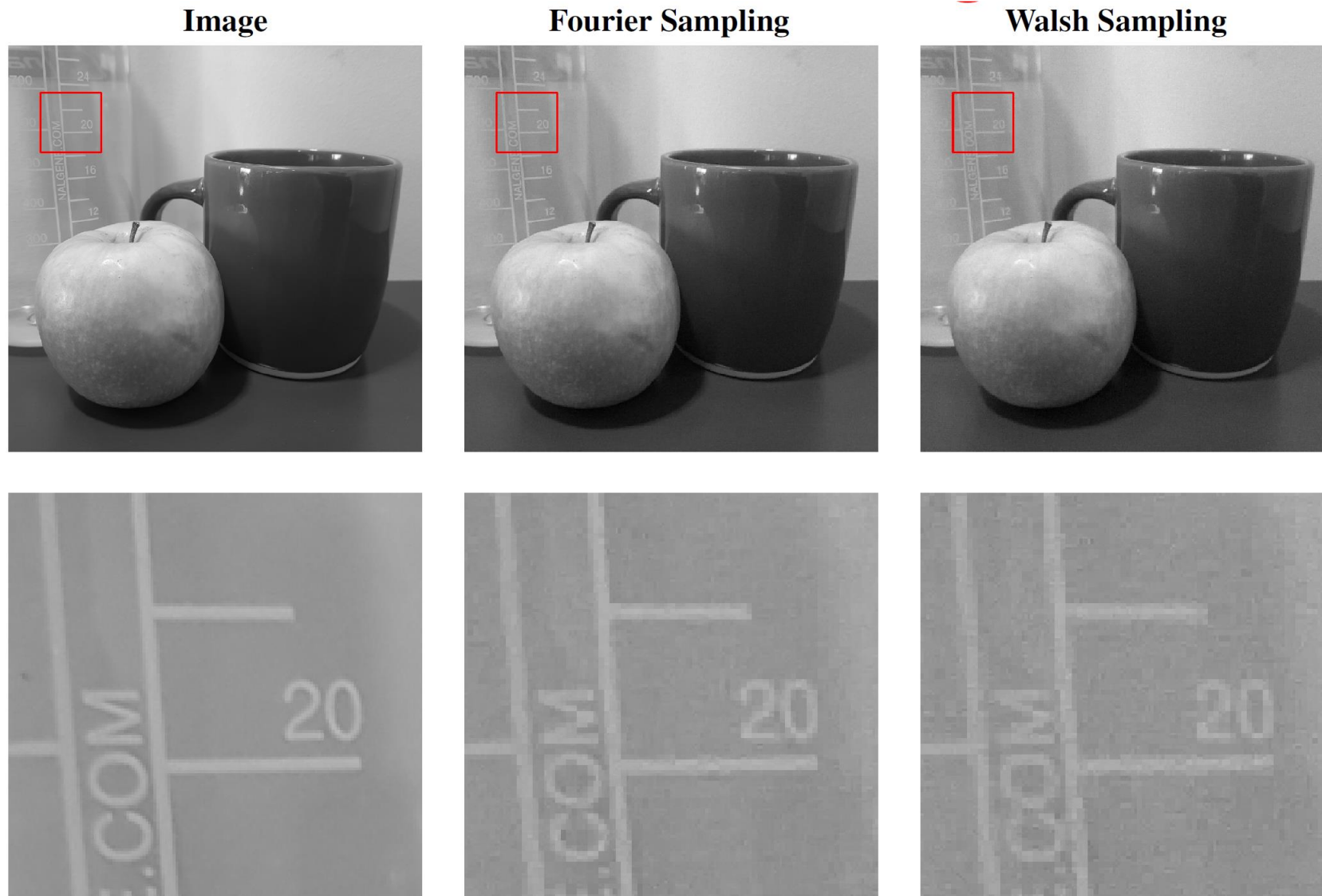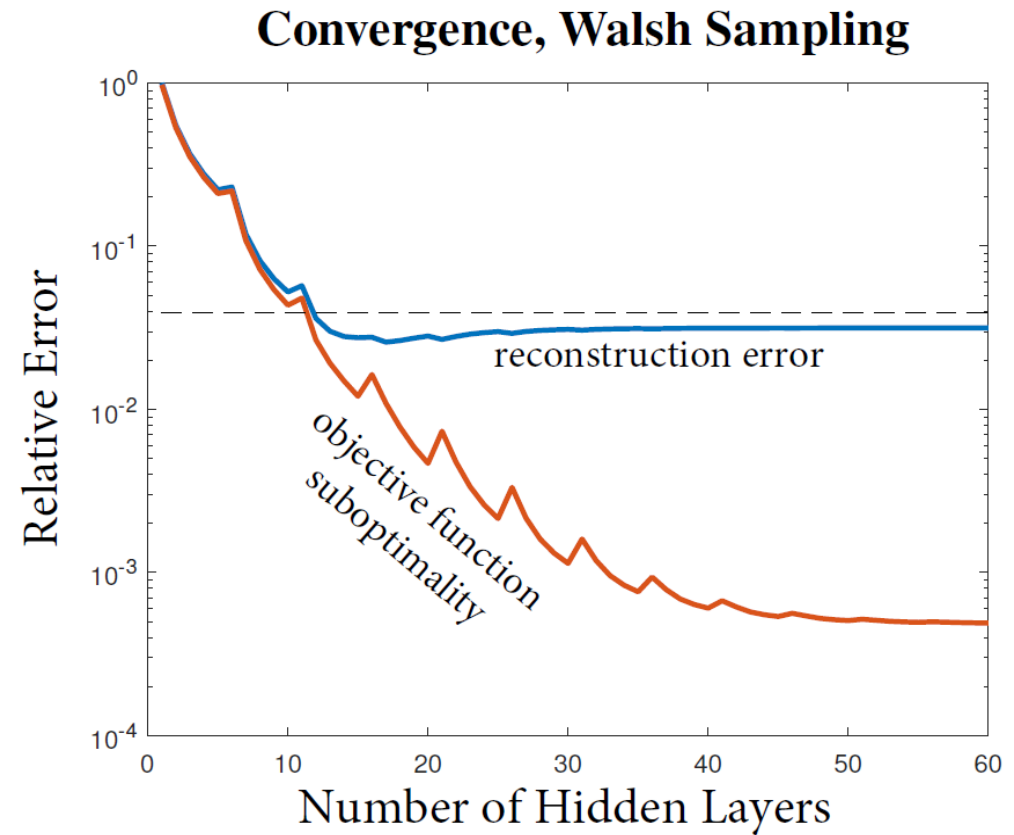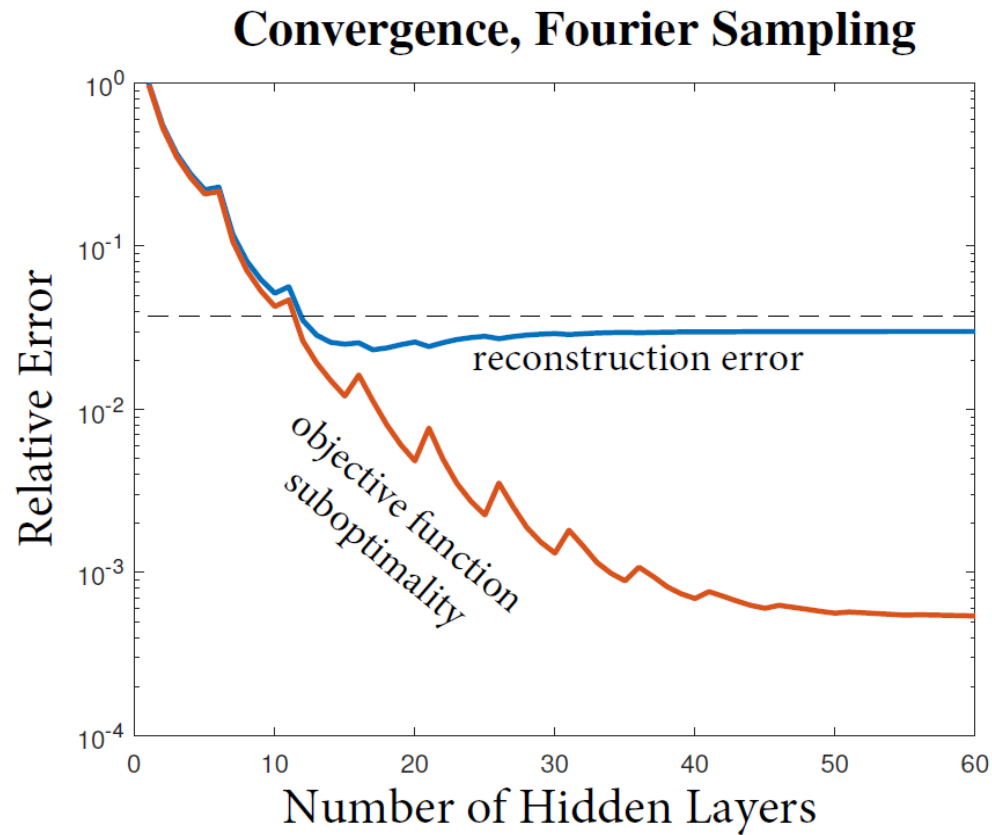**Image**  **Fourier Sampling**  **Walsh Sampling**



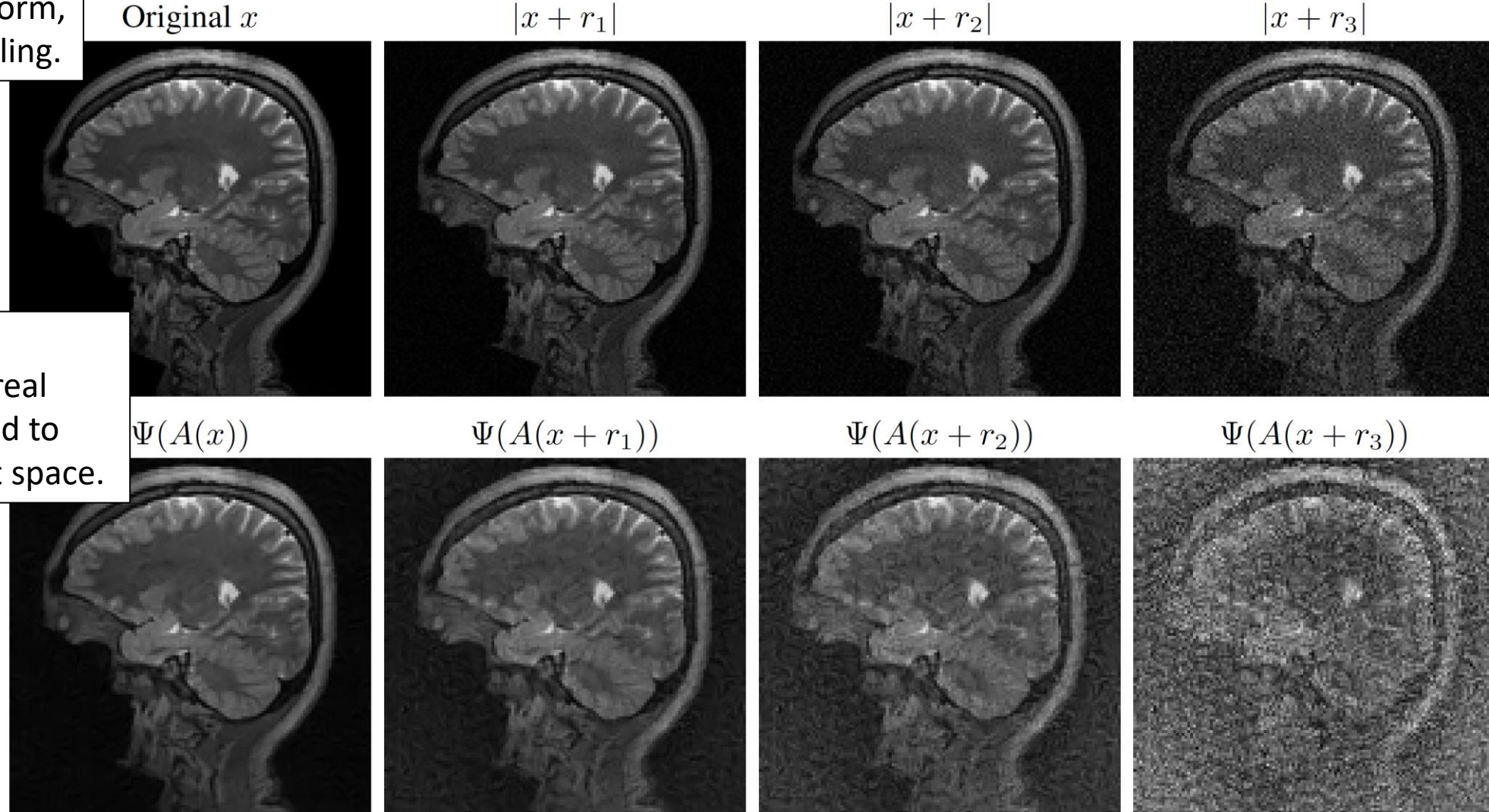Figure: Images corrupted with 2% Gaussian noise and reconstructed using 15% sampling.

# Demonstration of convergence
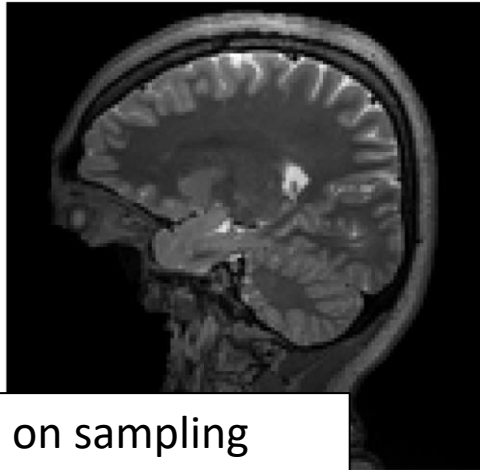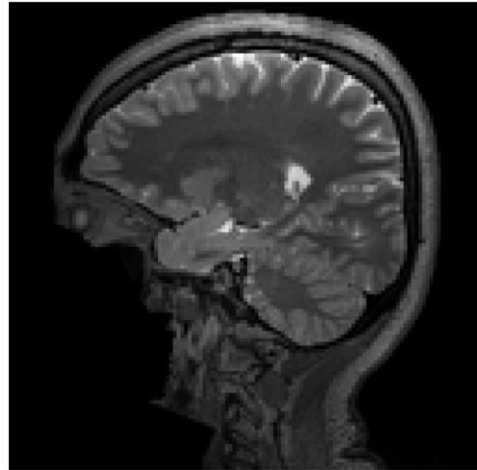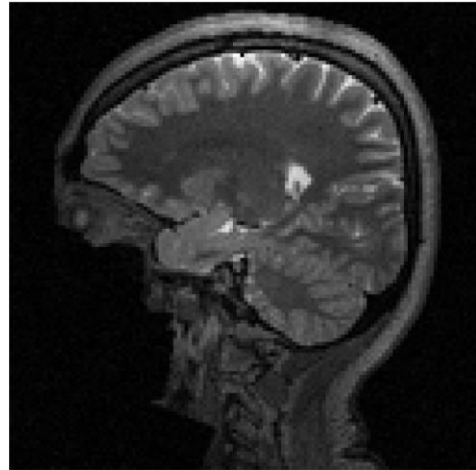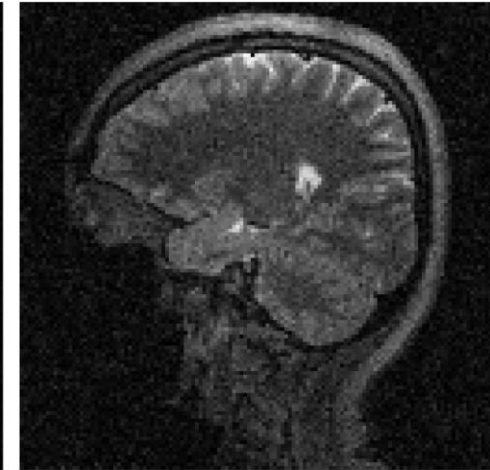
# Example of severe instability

MRI: discrete 2D Fourier transform, 60% subsampling.

Perturbations computed in real space, mapped to measurement space.



Original $x$  |  $|x + r_1|$  |  $|x + r_2|$  |  $|x + r_3|$

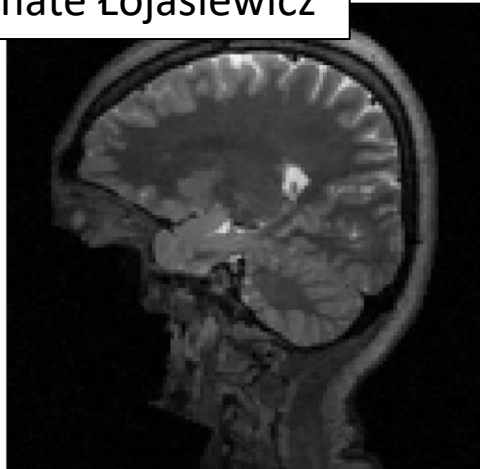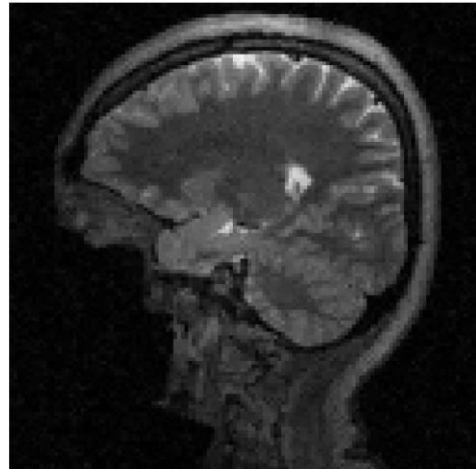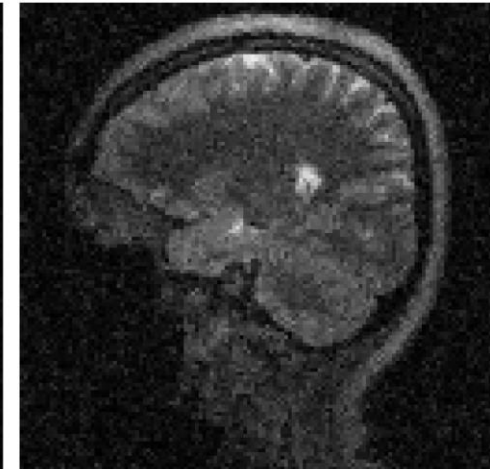$\Psi(A(x))$  |  $\Psi(A(x + r_1))$  |  $\Psi(A(x + r_2))$  |  $\Psi(A(x + r_3))$

- Zhu et al., "*Image reconstruction by domain-transform manifold learning,*" **Nature,** 2018.
- Antun et al., "*On instabilities of deep learning in image reconstruction and the potential costs of AI,*" **PNAS**, 2020.

# FIRENET: <u>provably</u> stable (even to adversarial examples) and accurate



Assumptions on sampling and approximate sparseness give approximate Łojasiewicz

• C., Antun, Hansen, "*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem*," **PNAS**, 2022.

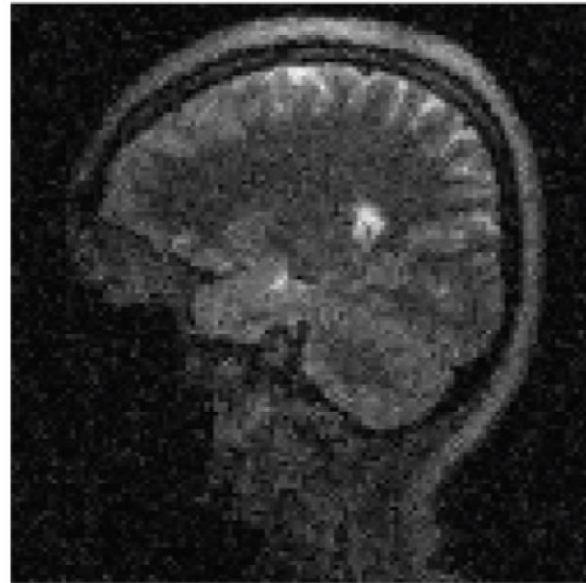# Stabilising unstable neural networks



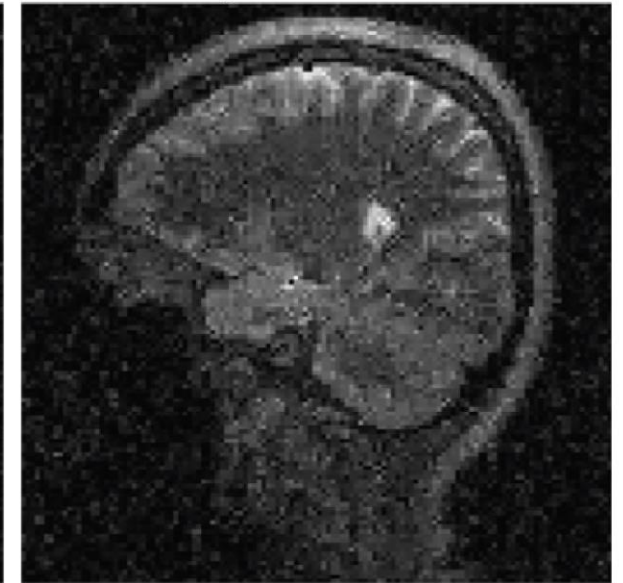$\Psi(\tilde{y}),\ \tilde{y} = Ax + e_3$     $\Phi\left(\tilde{y}, \Psi(\tilde{y})\right)$     FIRENET rec. from $y = Ax + \tilde{e}_3$     AUTOMAP+FIRENET rec. from $y = Ax + \hat{e}_3$

# Key pillars: stability and accuracy

MRI: discrete 2D Fourier transform, 15% subsampling.

All networks trained on 5000 images of ellipses



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

• C., Antun, 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: standard neural network architecture for imaging. Approx 4 million parameters.



• C., Antun, 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 with noise: stable but inaccurate
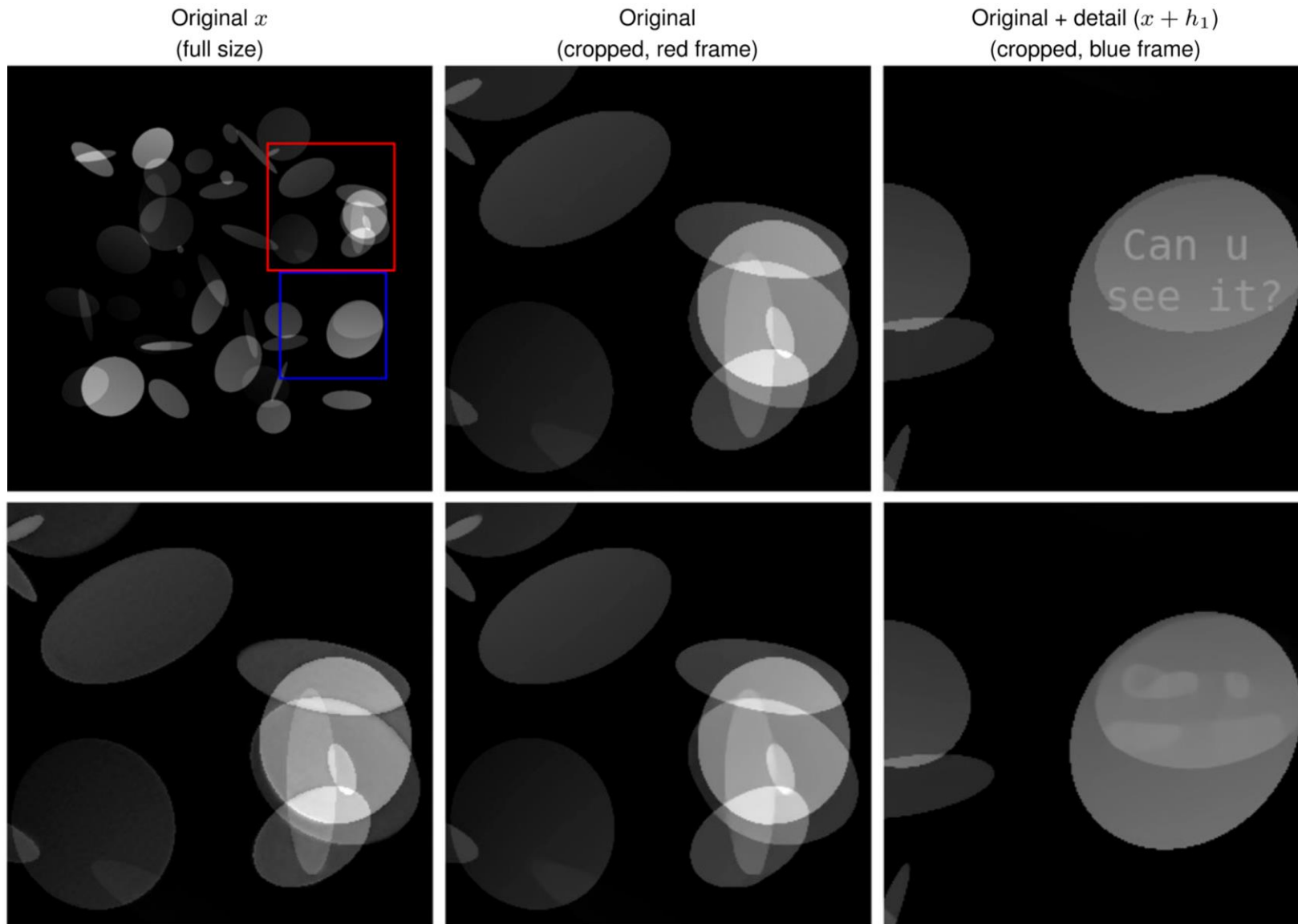


• C., Antun, Hansen, "*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem,*" **PNAS**, 2022.

# FIRENET: balances stability and accuracy?

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

Original $x$ (full size)

Original (cropped, red frame)

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

**Open problem:** use the toolkit to precisely prove theorems about *optimal* trade-offs.
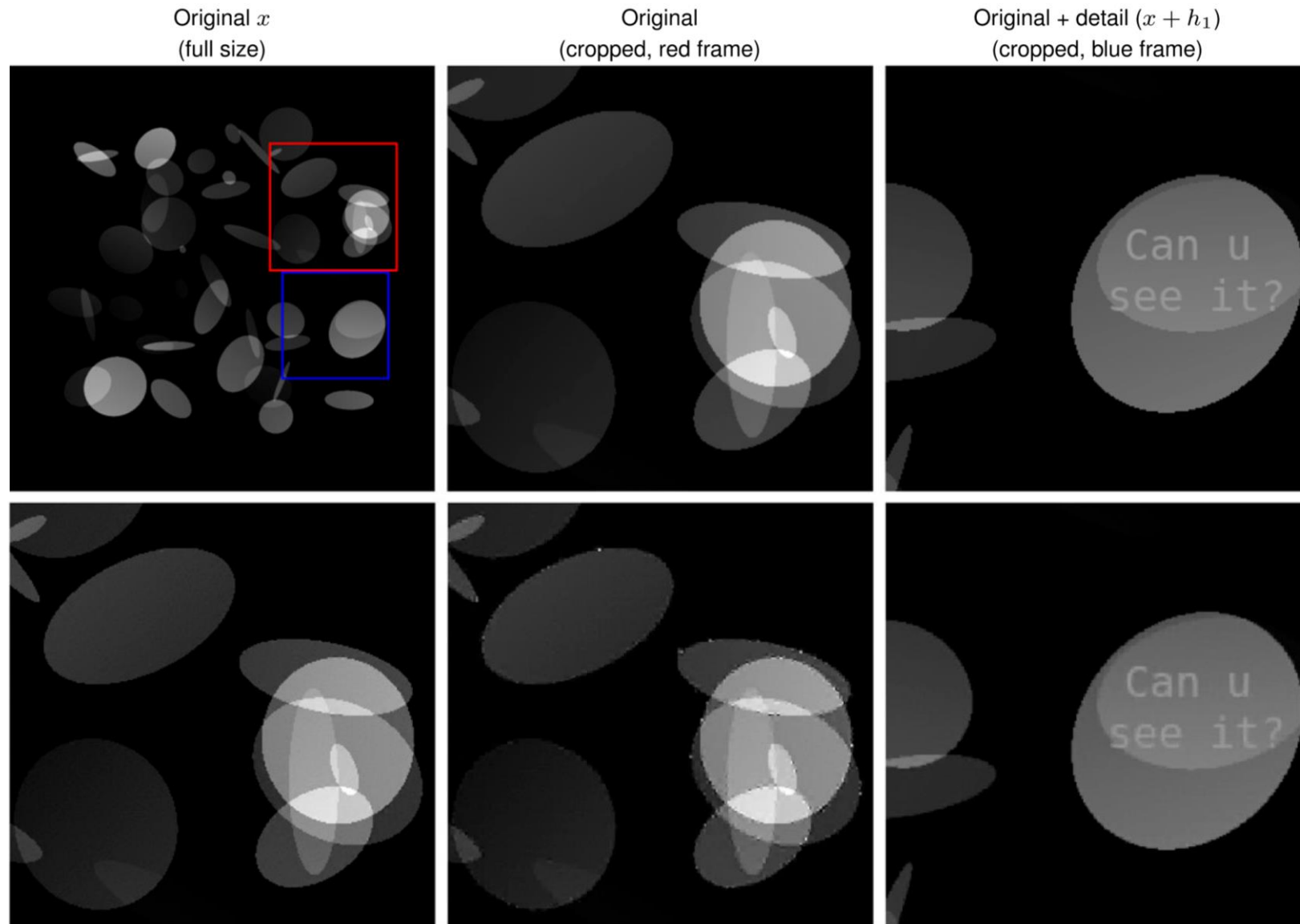
• C., Antun, Hansen, "*The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem,*" **PNAS**, 2022.
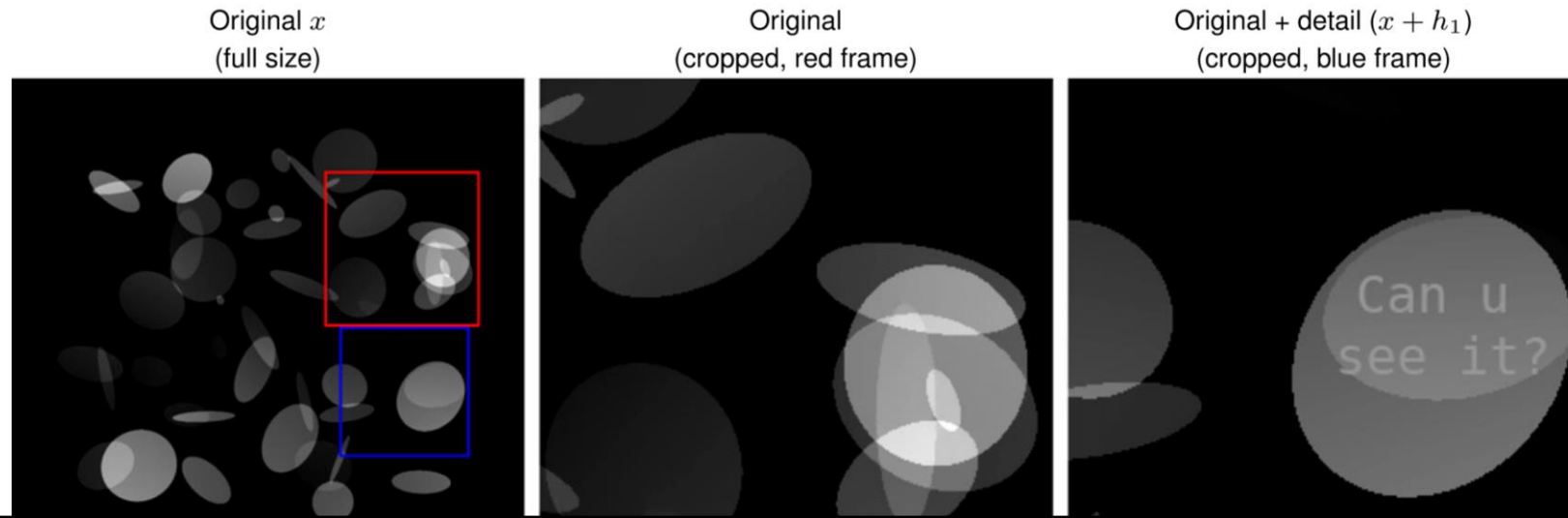
# Summary

## Need for foundations in AI/deep learning!

- **Paradox:** Nice linear inverse problems where stable and accurate neural network exists but cannot be trained!

- Trainability depends on
    - Accuracy desired.
    - Amount of training data.

- Specific conditions $\Rightarrow$ FIRENETs exp. convergence
                                    + withstand adversarial attacks.

- Trade-off between stability and accuracy in deep learning.