

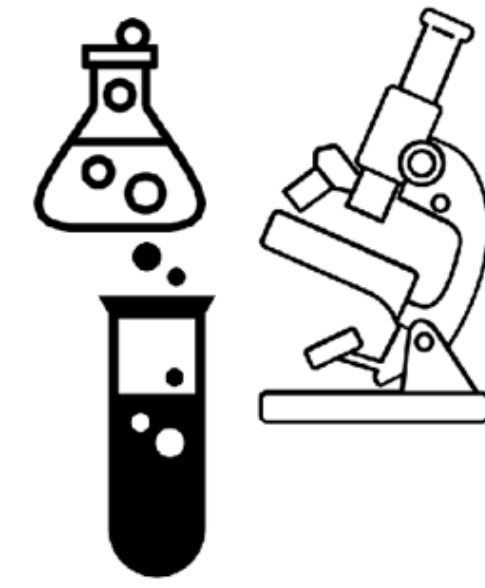
# A Learning Theoretic Perspective on Local Explainability

JEFFREY LI<sup>\*1</sup> VAISHNAVH NAGARAJAN<sup>\*2</sup> GREGORY PLUMB<sup>2</sup> AMEET TALWALKAR<sup>2,3</sup>

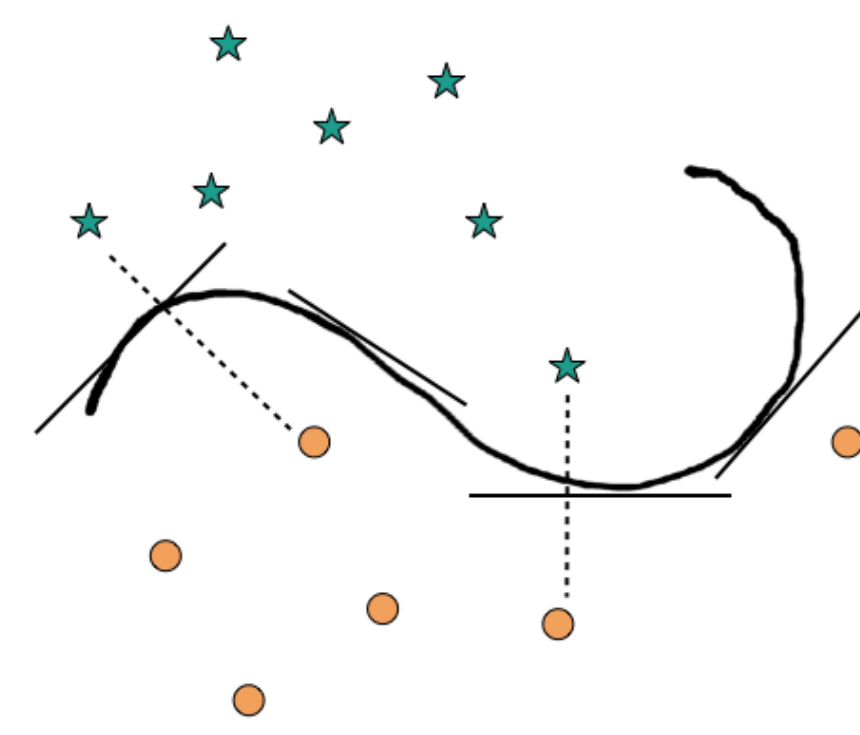
<sup>1</sup>University of Washington <sup>2</sup>Carnegie Mellon University <sup>3</sup>Determined AI

## OUR CONTRIBUTIONS

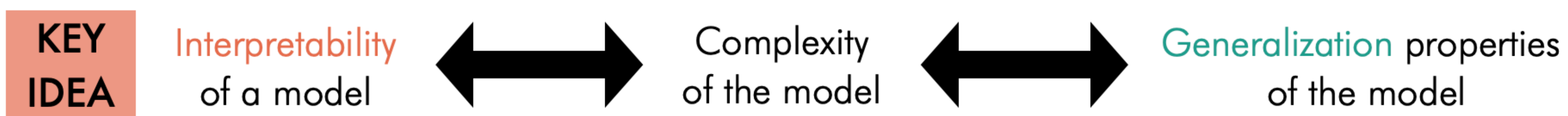
Interpretability has been a largely empirical field.



We establish one of the first connections between **interpretability** and **learning theory**:



1. We derive a generalization bound on test performance of a model in terms of its local explainability.
2. We address a new question: how well does the “quality” of local explanations generalize?



## PERFORMANCE GENERALIZATION

**Motivation:** Bounds based on complexity( $\mathcal{F}$ ) may not capture how “simple”  $f$  is!

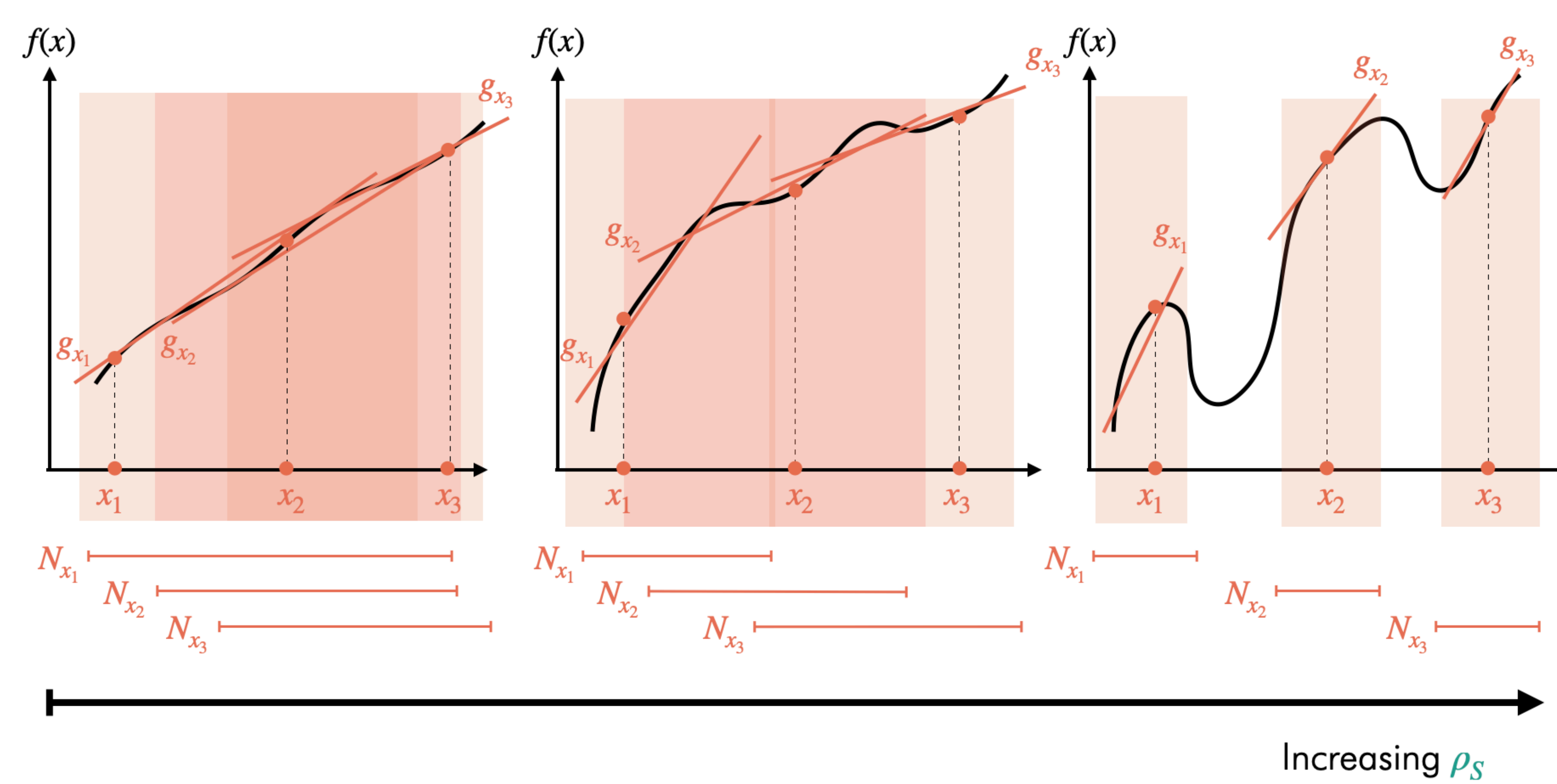
**Theorem 1:**

$$\mathbb{E}_D[(f(x) - y)^2] \leq \hat{\mathbb{E}}_S[(f(x) - y)^2] + \mathbb{E}_{x \sim D} \mathbb{E}_{x' \sim N_x} [(g_x(x) - f(x))^2] + \rho_S \cdot \mathcal{R}(\mathcal{G}_{local})$$

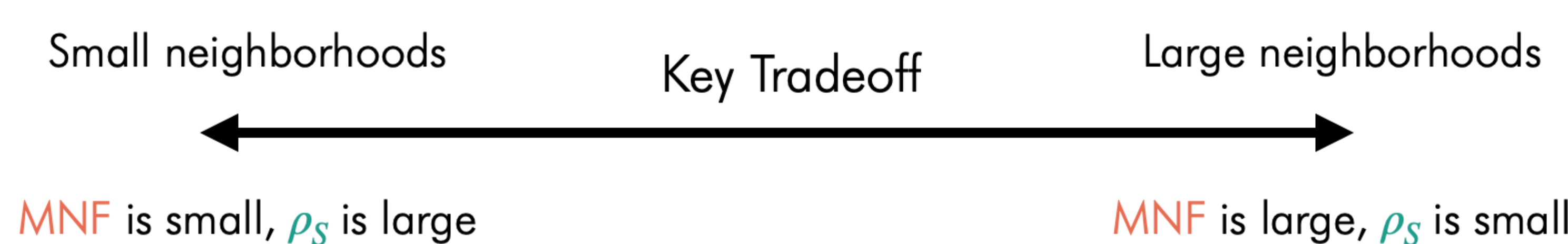
Test loss      Train loss (on dataset  $S$ )      Explanation quality (MNF)      Complexity of the system of local explanations

$$\mathcal{R}(\mathcal{G}_{local}) = \text{Rademacher complexity of } \mathcal{G}_{local} \approx O\left(\frac{1}{\sqrt{|S|}}\right) \text{ How complex each local explanation is}$$

$$\rho_S = \int_{x \in \mathcal{X}} \sqrt{\frac{1}{|S|} \sum_{x \in S} (p_{N_x}(x'))^2 dx'} \in [1, \sqrt{|S|}] \text{ How disjoint the neighborhoods of the training points are}$$



**Takeaway:** Better generalization when it is easier to locally approximate  $f$  on larger neighborhoods.



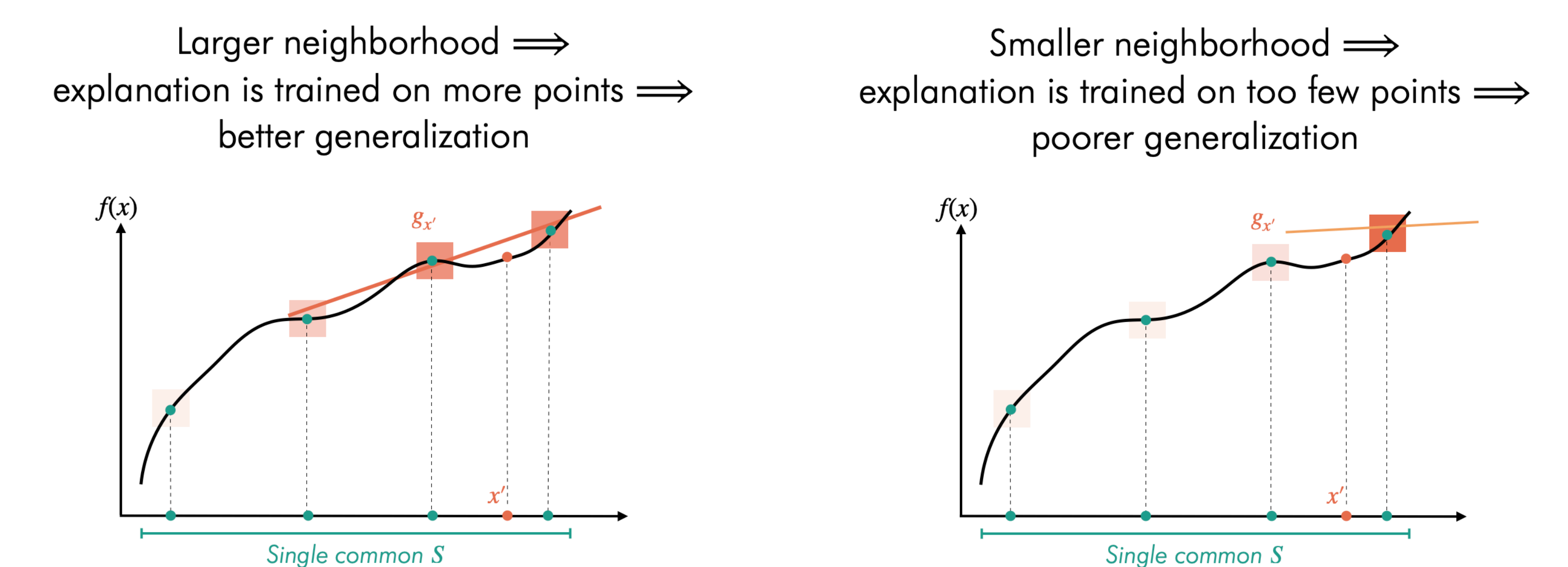
## EXPLANATION GENERALIZATION: RESULT

**Theorem 2:**

$$\mathbb{E}_{x \sim D} \mathbb{E}_{x' \sim N_x} [(g_x(x) - f(x))^2] \leq \hat{\mathbb{E}}_{x \sim S} \mathbb{E}_{x' \sim N_x} [(g_x(x) - f(x))^2] + \rho_S \cdot \mathcal{R}(\mathcal{G}_{local})$$

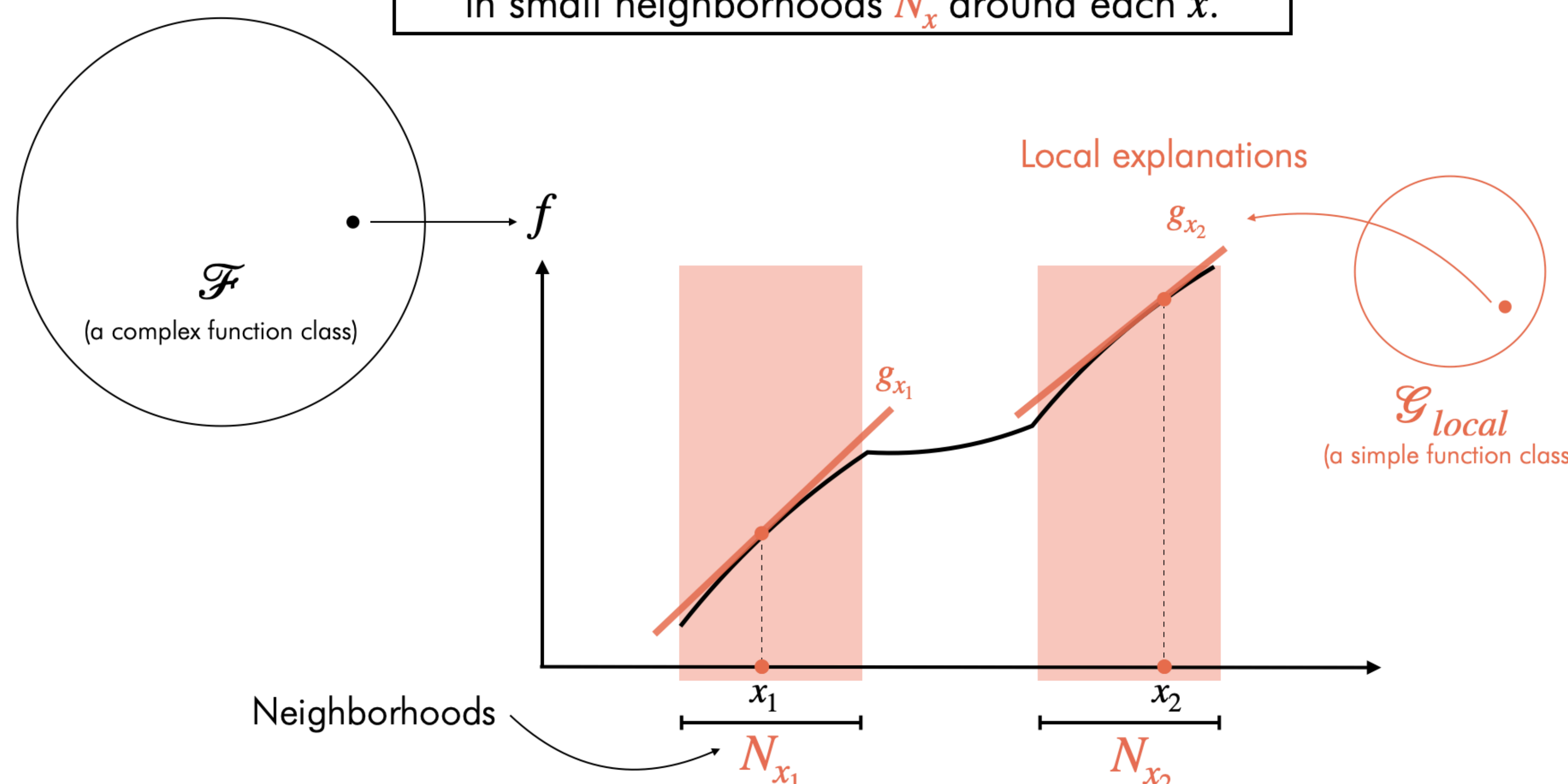
Test explanation quality (Test MNF)      Train explanation quality (Train MNF)      Complexity of the system of explanations

**Takeaway:** Better generalization when explanations can nicely fit training data that fall in a larger neighborhood.



## BACKGROUND: LOCAL EXPLANATIONS

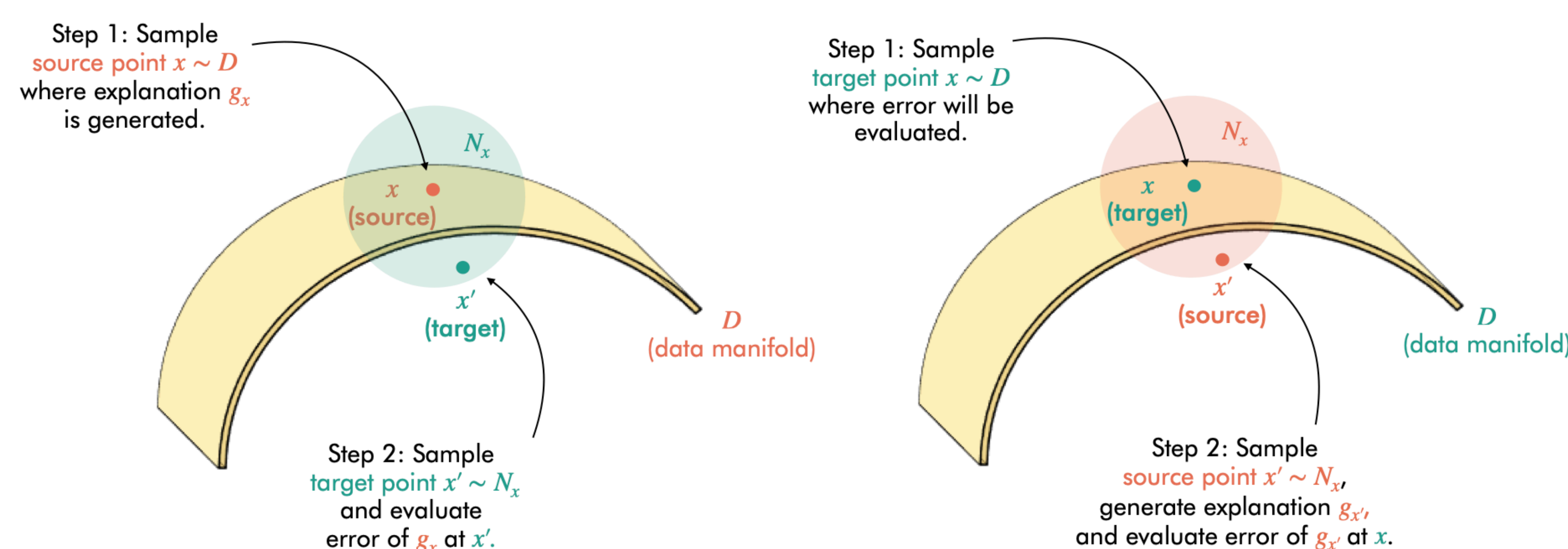
Approximate a function  $f$  via simple functions  $g_x$  in small neighborhoods  $N_x$  around each  $x$ .



## LOCAL EXPLANATION QUALITY

Standard measure: Neighborhood Fidelity (NF)  
 $\mathbb{E}_{x \sim D} \mathbb{E}_{x' \sim N_x} [(g_x(x') - f(x'))^2]$

We propose: Mirrored Neighborhood Fidelity (MNF)  
 $\mathbb{E}_{x \sim D} \mathbb{E}_{x' \sim N_x} [(g_x(x) - f(x))^2]$

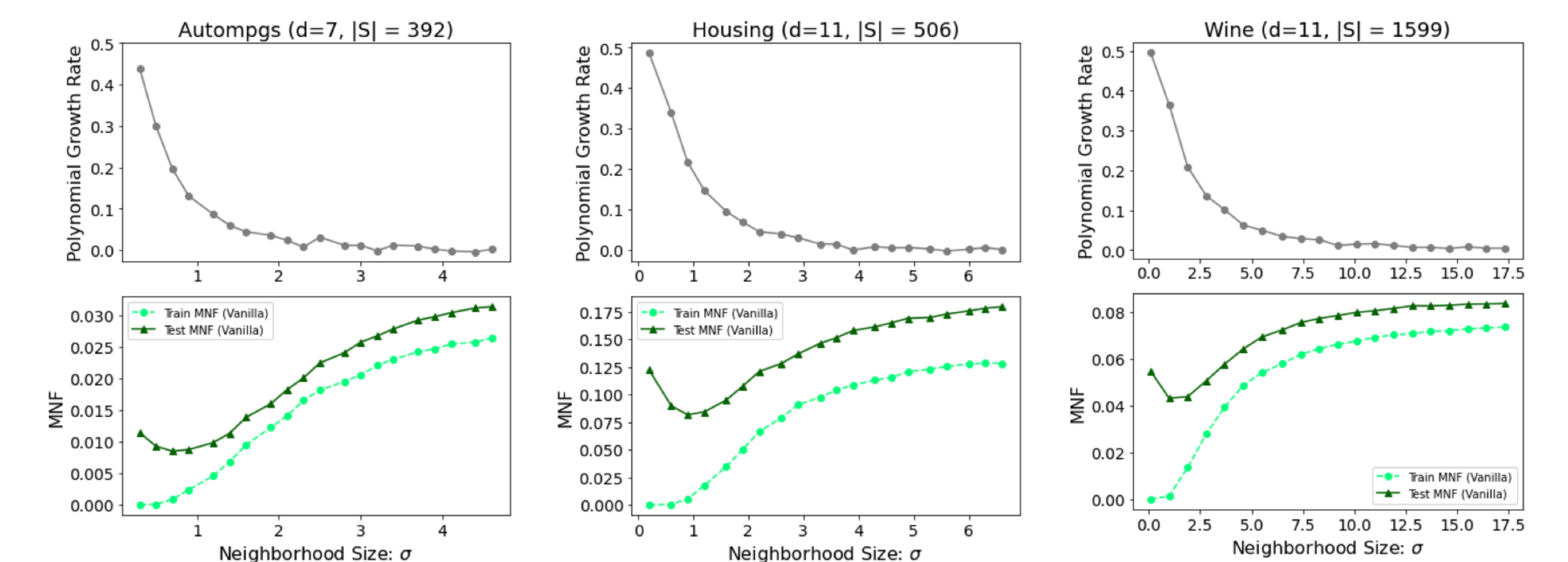


Unlike NF, MNF does not evaluate fit of explanations on off-manifold data!

- ⇒ more amenable to theoretical analysis
- ⇒ more robust to irregular off-manifold behavior of  $f$ .

## EXPERIMENTS

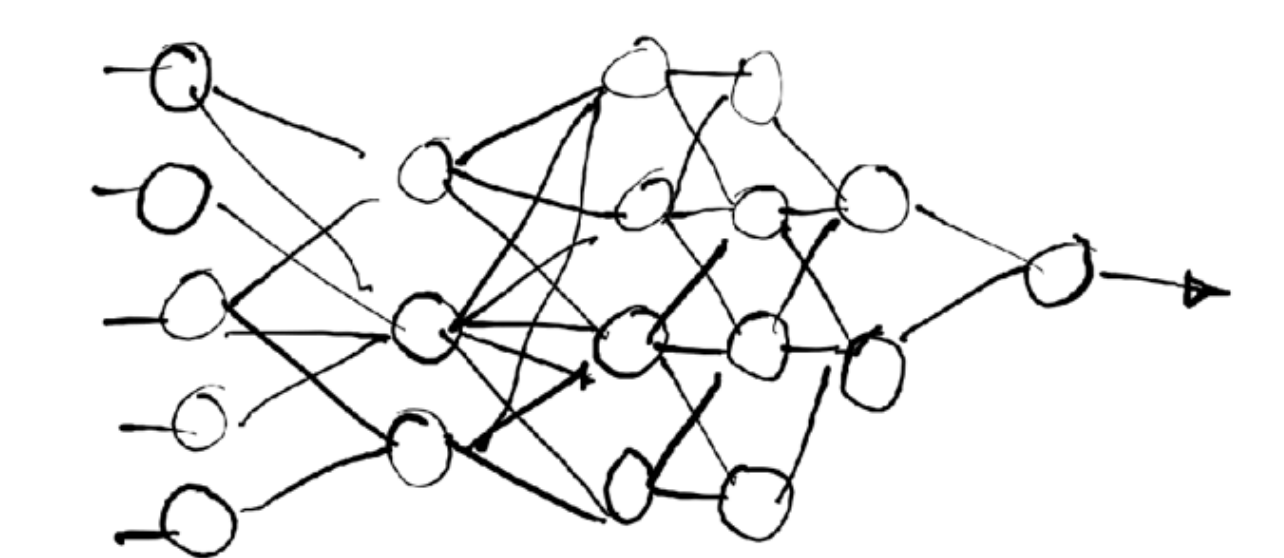
- Are there neighborhood widths s.t.  $\rho_S = o(|S|^{0.5})$  while Train MNF is small? **Yes!**
- Do wider neighborhoods bring the generalization gap down? **Yes!**



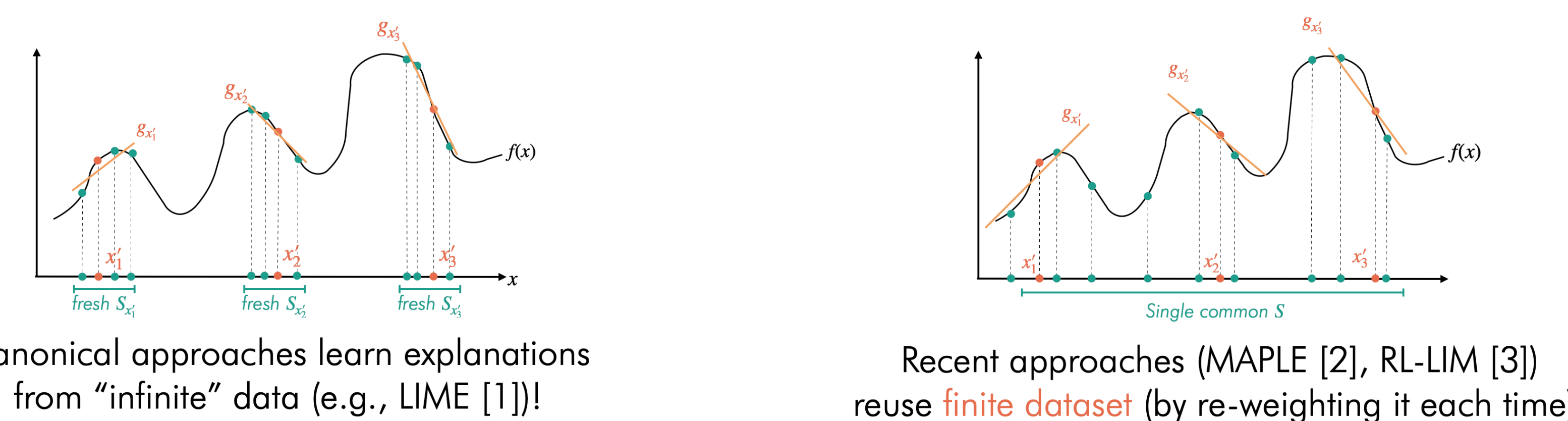
Exponent of  $\rho_S$  (top) and test & train MNF (bottom) on UCI datasets

## FUTURE WORK

1. Extend bounds to high-dimensional datasets.
2. Explore these bounds for NNs.
3. When is MNF better than NF in practice?



## EXPLANATION GENERALIZATION: MOTIVATION



Finite-sample-based approaches could potentially “overfit their explanations”!

What determines the quality of these explanations on unseen data?

## REFERENCES

1. “Why should I trust you?: Explaining the predictions of any classifier.” Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, ACM SIGKDD, 2016.
2. “Model Agnostic Supervised Local Explanations”, Gregory Plumb, Denali Molitor and Ameet S. Talwalkar, NeurIPS 2018
3. “RL-LIM: Reinforcement learning-based locally interpretable modeling”, Jinsung Yoon, Sercan O. Arık, and Tomas Pfister, 2019
4. “Uniform convergence may be unable to explain generalization in deep learning”. Vaishnavh Nagarajan and J. Zico Kolter, NeurIPS 2019
5. “Understanding deep learning requires rethinking generalization”, Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, Oriol Vinyals, ICLR’ 17