

# shapiq: Shapley Interactions for Machine Learning

Maximilian Muschalik<sup>1</sup>, Hubert Baniecki<sup>2</sup>, Fabian Fumagalli<sup>3</sup>, Patrick Kolpaczki<sup>4</sup>, Barbara Hammer<sup>3</sup>, and Eyke Hüllermeier<sup>1</sup>

How do I measure **interactions** between multiple features for **black box** models beyond feature attributions?



I want to use Shapley values for **other ML applications**. How do I compute them?

## Explain Models with Shapley Interactions

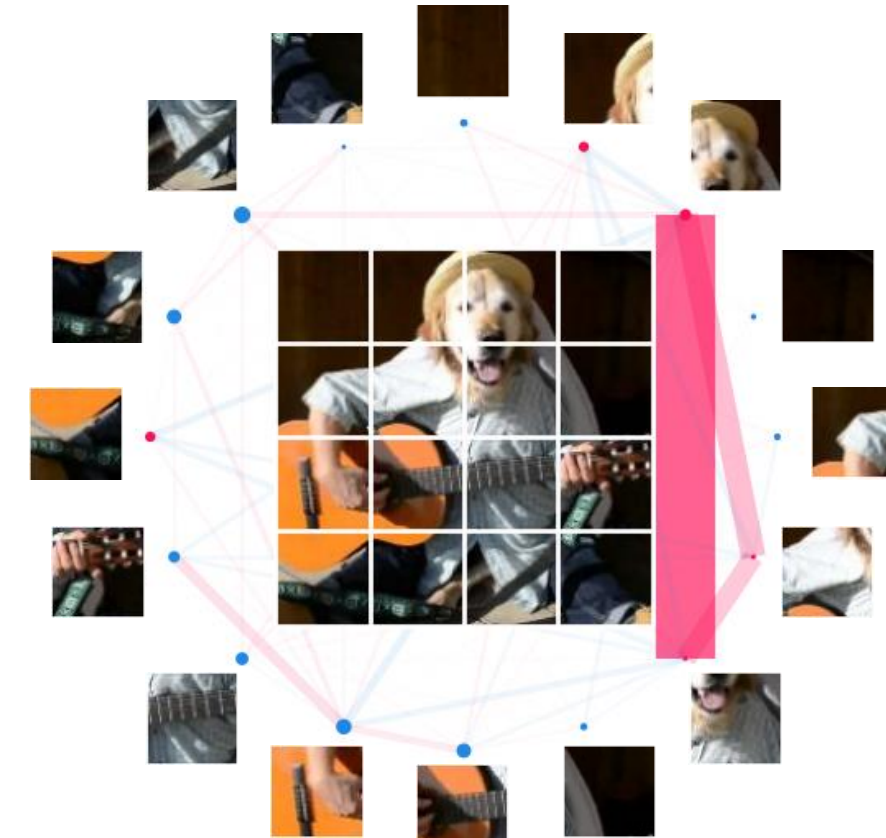
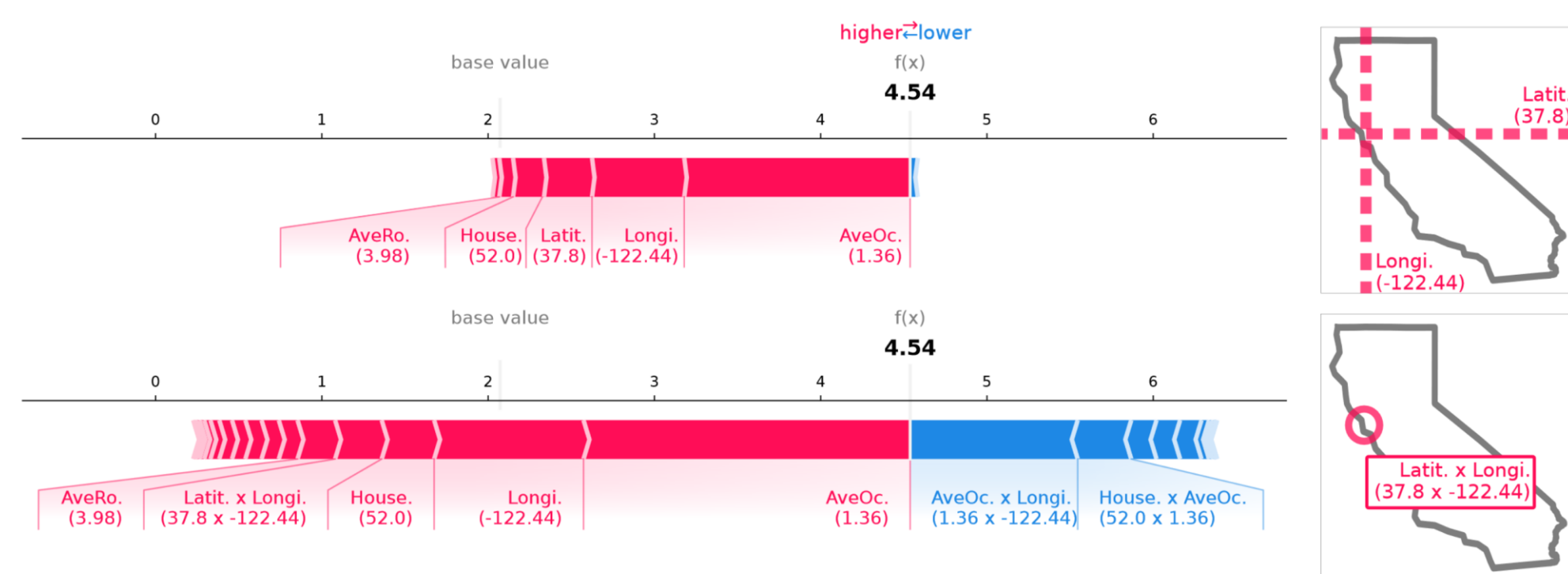
Explaining models with shapiq is **easy**:

- Agnostic Explainer and Imputers
- Tree Explainer

```
# get your data and model
X, model = ...
from shapiq import Explainer
# create an explainer object
explainer = Explainer(model=model, data=X, max_order=2)
# get the feature interactions for the first observation
interaction_values = explainer.explain(X[0], budget=1024)
# visualize the 2-order feature interactions
interaction_values.force_plot(feature_names=...)
```

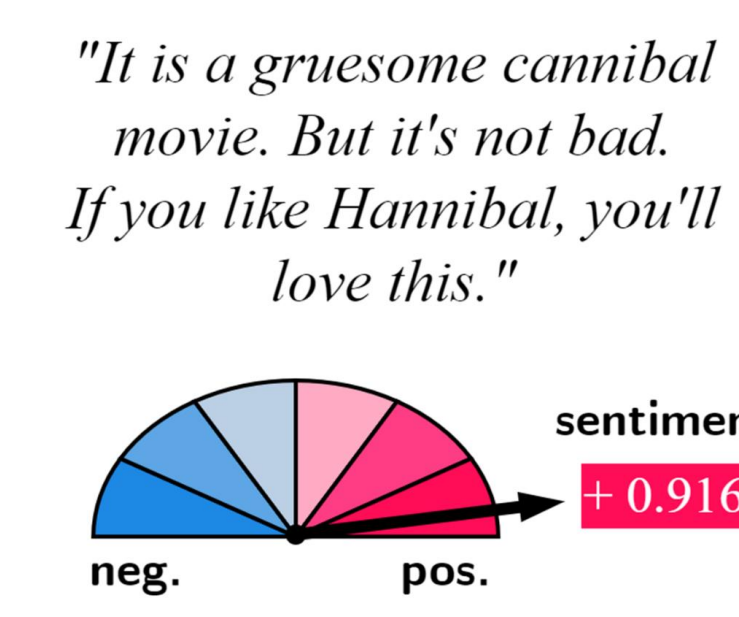
“Does the **location** of my property affect its price?”

“Why is this a **dog**?”



“How does my **language model** predict a positive sentiment?”

Sentiment Analysis Model



Explanation

SHAP: It is a **gruesome** **cannibal** **movie**. **But** **it's** **not** **bad**. **If** **you** **like** **Hannibal**, **you'll** **love** **this**.

SHAP-IQ: It is a **gruesome** **cannibal** **movie**. **But** **it's** **not** **bad**. **If** **you** **like** **Hannibal**, **you'll** **love** **this**.

## Game Theory for General ML Applications

Any Model (e.g., torch, sklearn, ...)

Tree Model (e.g., xgboost, lightgbm, ...)

Any Value Function (as a callable)

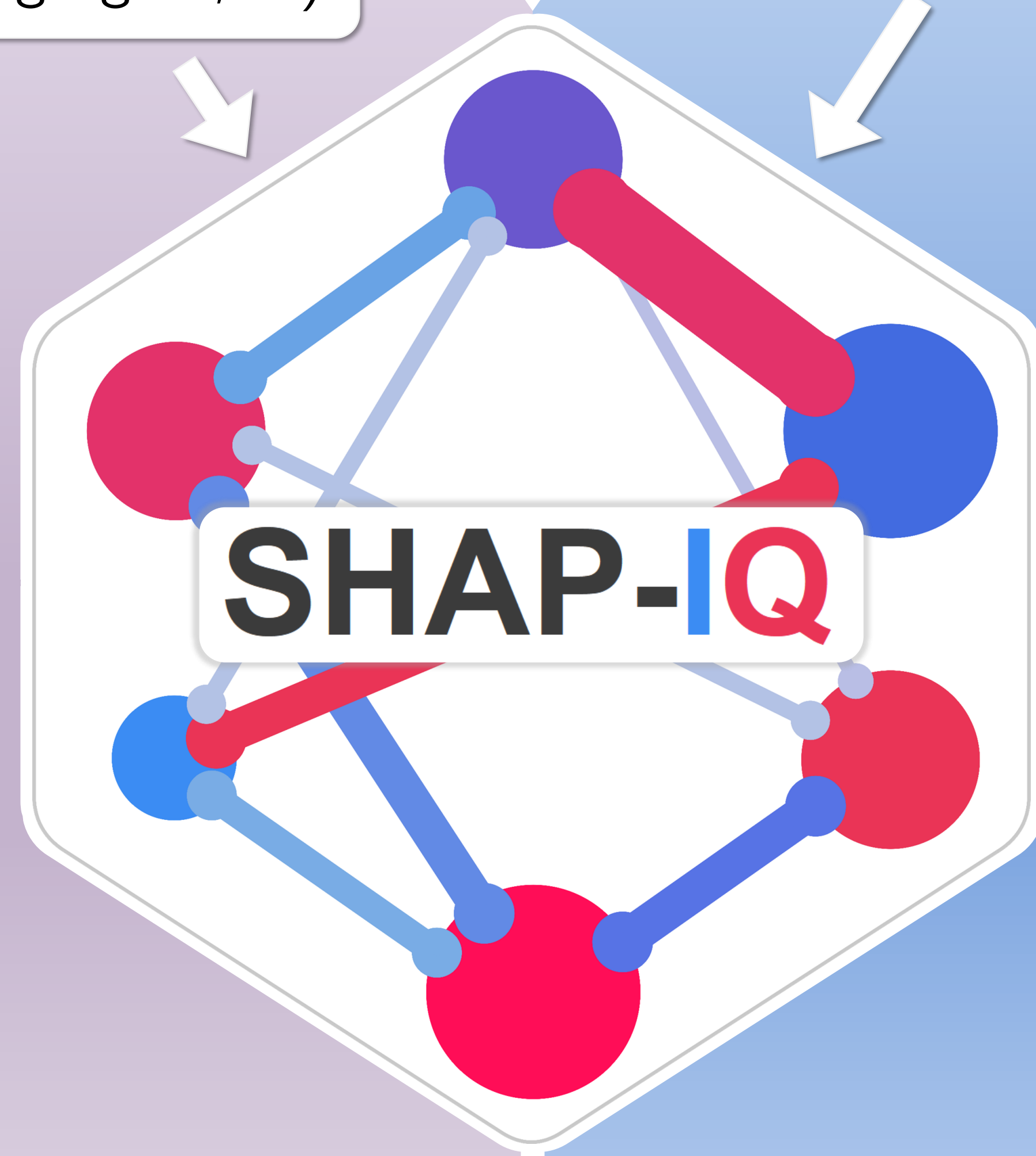
$v : \mathcal{P}(N) \rightarrow \mathbb{R}$

shapiq includes:

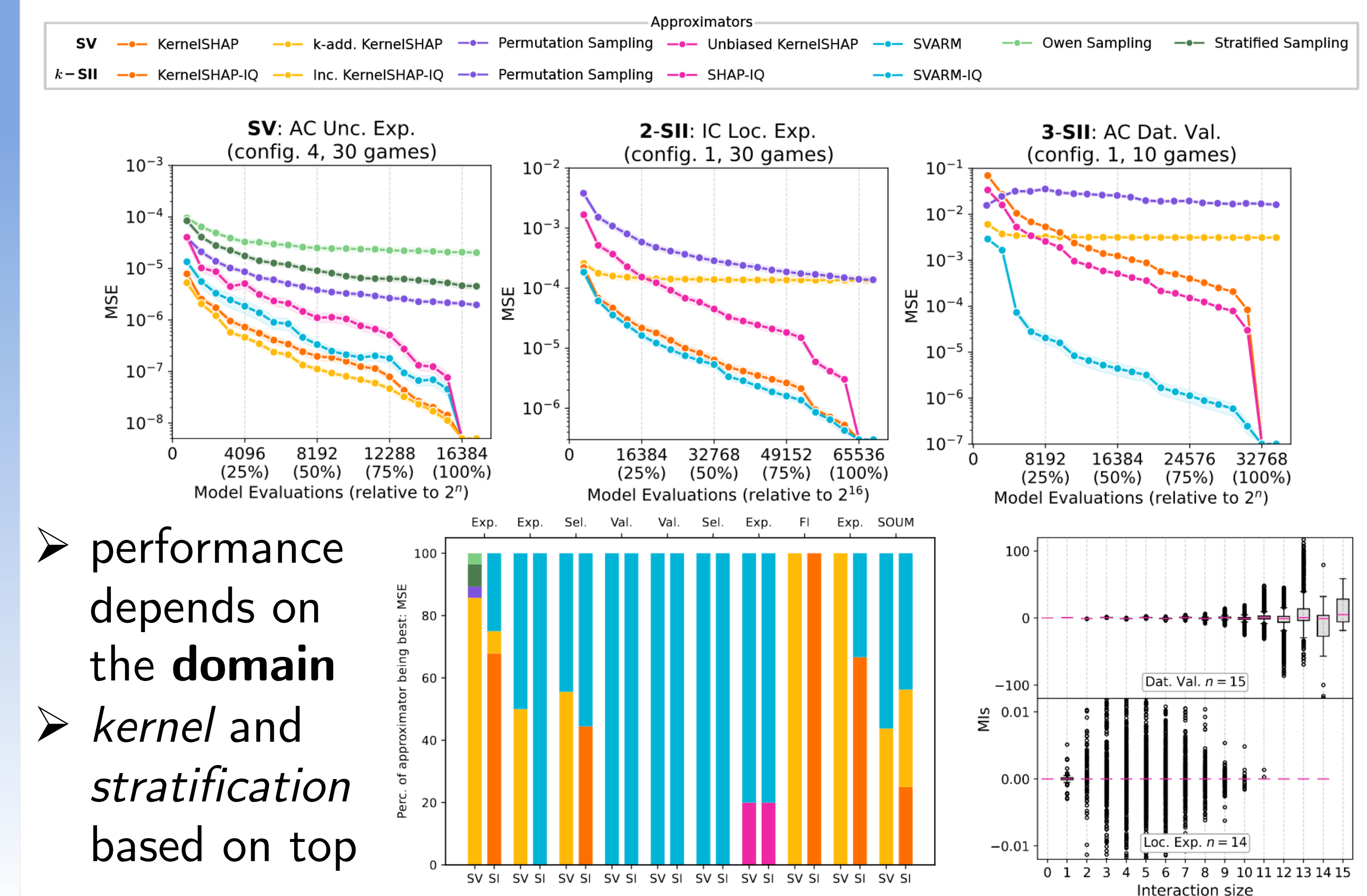
- 20 concepts (Shapley value and interactions, Banzhaf value and interactions, Faithful Shapley, Generalized values, Möbius, Core, ...)
- 14 state-of-the-art **approximators** and **exact computers**

```
import shapiq
class CountGame(shapiq.Game):
def __init__(self, n_players): ...
def value_function(coalitions):
# define the worth of a coalition
return np.sum(coalitions, axis=1)
game = CountGame(n_players=12)
approx = shapiq.KernelSHAPIQ(n=12)
si = approx(game=game, budget=1000)
# compute the Moebius transform exactly
exact = shapiq.ExactComputer(game, 12)
mi = exact(index='Moebius')
print(si[(3, 7)], mi[(3,)]) # get values
```

Class	Shapley Interactions	Shapley Values
Approximator	KernelSHAP-IQ	KernelSHAP
	Inconsistent KernelSHAP-IQ	k <sub>ADD</sub> -SHAP
	Faith-SHAP	Owen Sampling
	SHAP-IQ	Unbiased KernelSHAP
	SVARM-IQ	SVARM
	Permutation Sampling (SII)	Permutation Sampling (SV)
Computer	Permutation Sampling (STII)	Stratified Sampling
		Möbius Converter
		Exact Computer



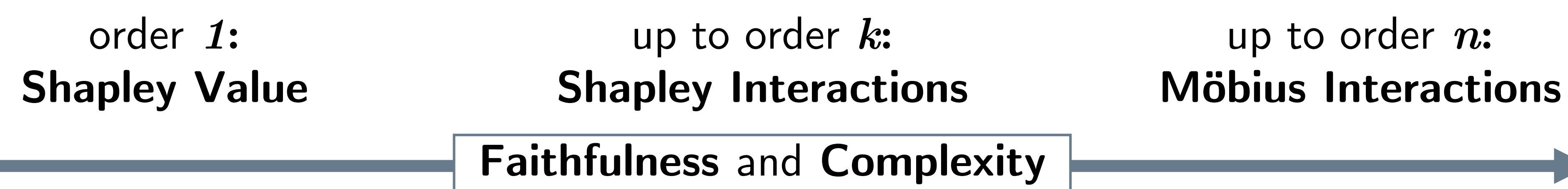
## Evaluation of Approximators on the Benchmark



**Benchmark of 11 ML domains** (e.g., explanation, data valuation, uncertainty quantification, ...)

**Games: 100** benchmark games with more than **2000** pre-computed configurations

**Interpretation:** Shapley interactions generalize the Shapley value beyond individual effects up to **any-order** and capture **synergies** between features.



- performance depends on the **domain**
- **kernel** and **stratification** based on top