

SVARM-IQ: Efficient Approximation of Any-order Shapley Interactions through Stratification

Patrick Kolpaczki¹, Maximilian Muschalik², Fabian Fumagalli³, Barbara Hammer³, and Eyke Hüllermeier²



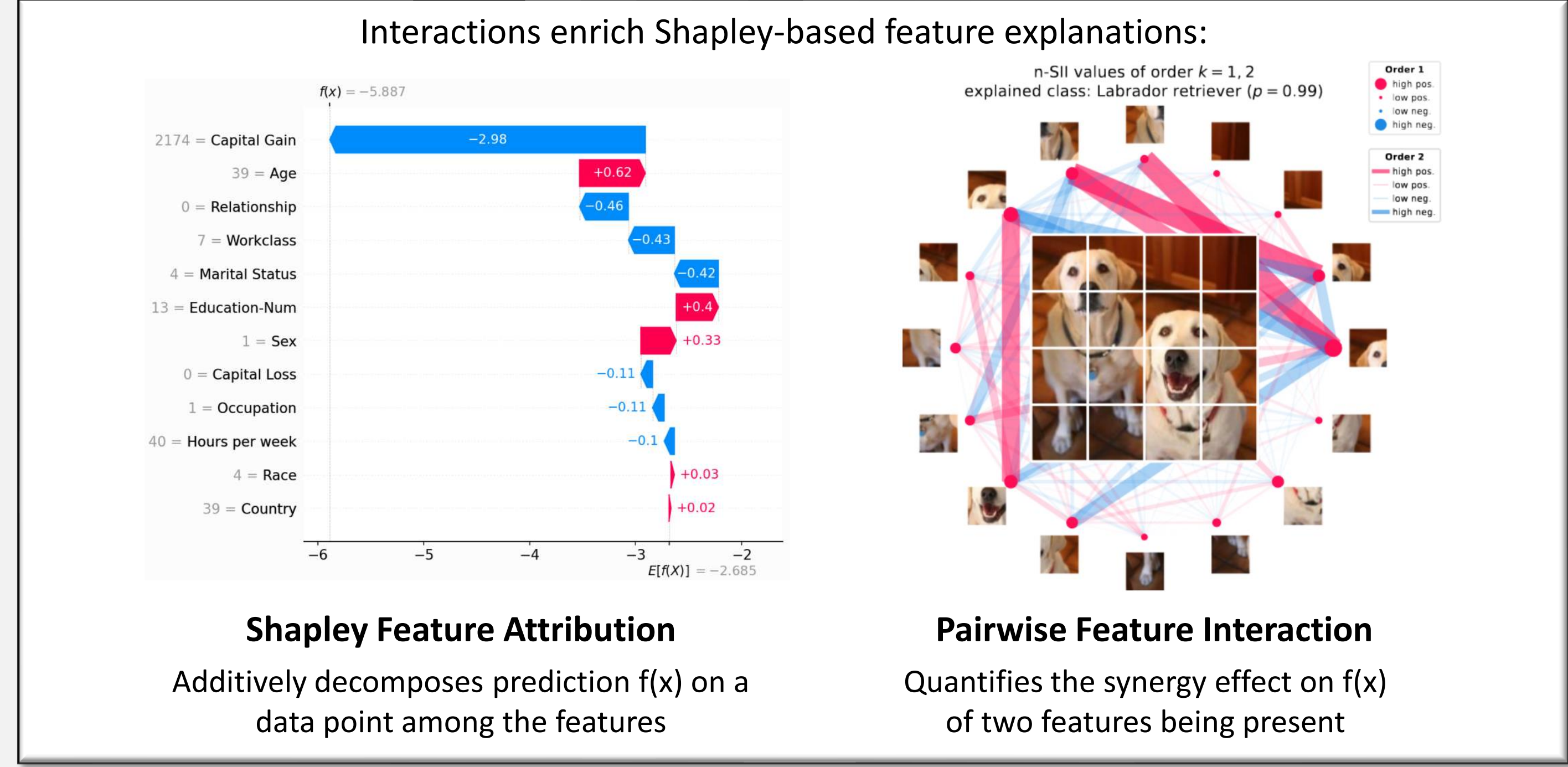
1) Paderborn University, Germany

2) LMU Munich, Germany

3) Bielefeld University, Germany



Motivation: Feature Explanations



Contribution

- Novel approximation algorithm for all Shapley-based Interactions:**
- Powerful combination of stratified representation + update mechanism
 - Applicable to **all orders** (pairs, triples, etc.) and **indices** (SII, STI, etc.) simultaneously
 - Novel theoretical guarantees & state-of-the-art empirical performance
 - ✓ **Model-agnostic** / domain-independent
→ *Applicable to any model and data, and even outside of explainability and ML*
 - ✓ No hyperparameters → *No fine-tuning*
 - ✓ Estimates available at any time → *Budget can be cut and extended arbitrarily*

Shapley-based Interaction Indices

- Player set** $\mathcal{N} = \{1, \dots, n\}$ → Features, datapoints, neurons, base learners etc.
 - Value function** $\nu: \mathcal{P}(\mathcal{N}) \rightarrow \mathbb{R}$ → Predicted value, generalization performance with $\nu(\emptyset) = 0$
- Definition: Shapley Value** (Shapley, 1953)
- $$\phi_i = \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \frac{1}{n \cdot \binom{n-1}{|S|}} \cdot [\nu(S \cup \{i\}) - \nu(S)]$$
- Marginal contribution $\Delta_i(S)$: Increase in collective benefit when i joins S .
- Unique solution to fulfill desirable axioms: Symmetry, Additivity, Null-Property, Efficiency
 - Computational effort scales **exponentially** with n : 2^n coalitions in total

- Definition: Cardinal Interaction Indices** (Fujimoto et al., 2006)
- For pairs i, j :
$$I_{i,j} = \sum_{S \subseteq \mathcal{N} \setminus \{i,j\}} \lambda_{2,|S|} \cdot [\nu(S \cup \{i,j\}) - \nu(S \cup \{i\}) - \nu(S \cup \{j\}) + \nu(S)]$$
- For subsets K of order k :
$$I_K = \sum_{S \subseteq \mathcal{N} \setminus K} \lambda_{k,|S|} \cdot \Delta_K(S)$$
 with $\Delta_K(S) := \sum_{W \subseteq K} (-1)^{|K|-|W|} \cdot \nu(S \cup W)$
- Discrete derivative $\Delta_K(S)$: Synergy effect of K at the presence of S .
- The weights $\lambda_{k,|S|}$ define the specific Interaction Index:
- Shapley Interaction Index (SII):** $\lambda_{k,|S|}^{\text{SII}} = \frac{1}{(n-k+1) \cdot \binom{n-k}{|S|}}$
 - Shapley-Taylor Interaction Index (STI):** $\lambda_{k,|S|}^{\text{STI}} = \frac{k}{n \cdot \binom{n-1}{|S|}}$
- And many more:
- Faithful-Shapley Interaction Index (FSI)
 - Banzhaf Interaction Index (BII)

- Fixed-budget approximation problem:**
- Given cooperative game (N, ν) with unknown Interaction scores I_K for all $K \subseteq N$ of order k
 - Budget B : Allowed number of evaluations of ν (bottleneck due to model access)
Model evaluations (inference, retraining) pose bottleneck on runtime rather than arithmetic operations
 - Minimize mean squared error (MSE) averaged over all estimates \hat{I}_K :
$$\frac{1}{\binom{n}{k}} \sum_{K \subseteq \mathcal{N}} \mathbb{E} [(\hat{I}_K - I_K)^2]$$

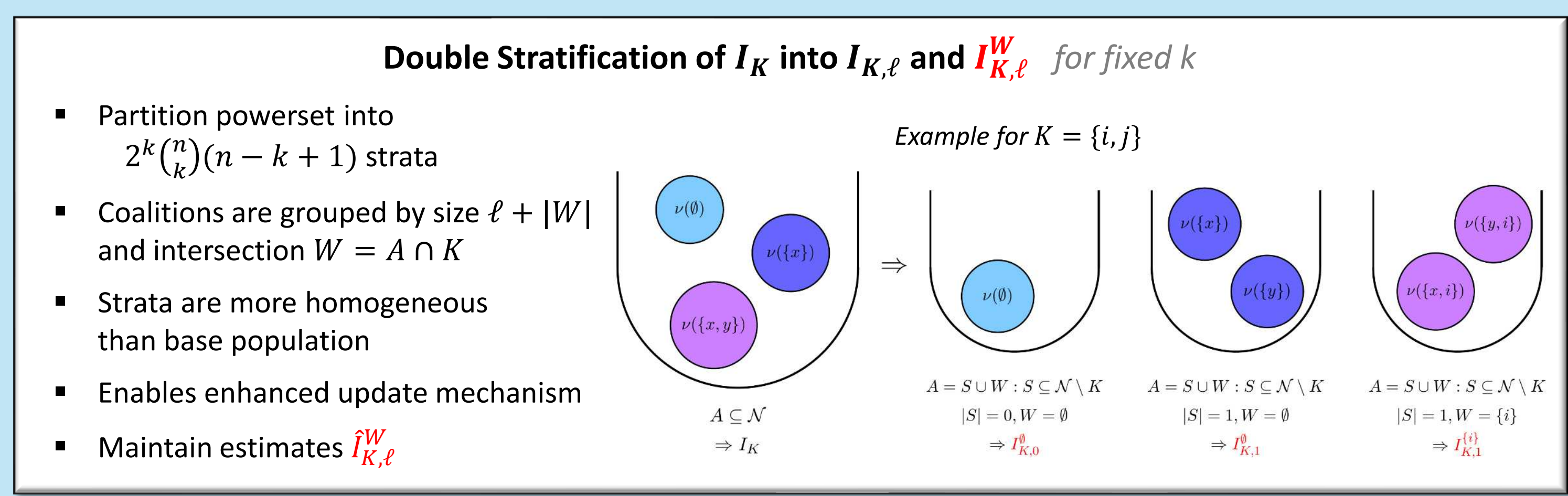
Approximation Algorithm: SVARM-IQ

Stratified Representation

$$I_K = \sum_{\ell=0}^{n-k} \binom{n-k}{\ell} \lambda_{k,\ell} \sum_{W \subseteq K} (-1)^{k-|W|} \cdot I_{K,\ell}^W$$

with
$$I_{K,\ell}^W := \frac{1}{\binom{n-k}{\ell}} \sum_{\substack{S \subseteq \mathcal{N} \setminus K \\ |S|=\ell}} \nu(S \cup W)$$

Stratum: Average worth of coalitions of size $\ell+|W|$, containing only W out of K



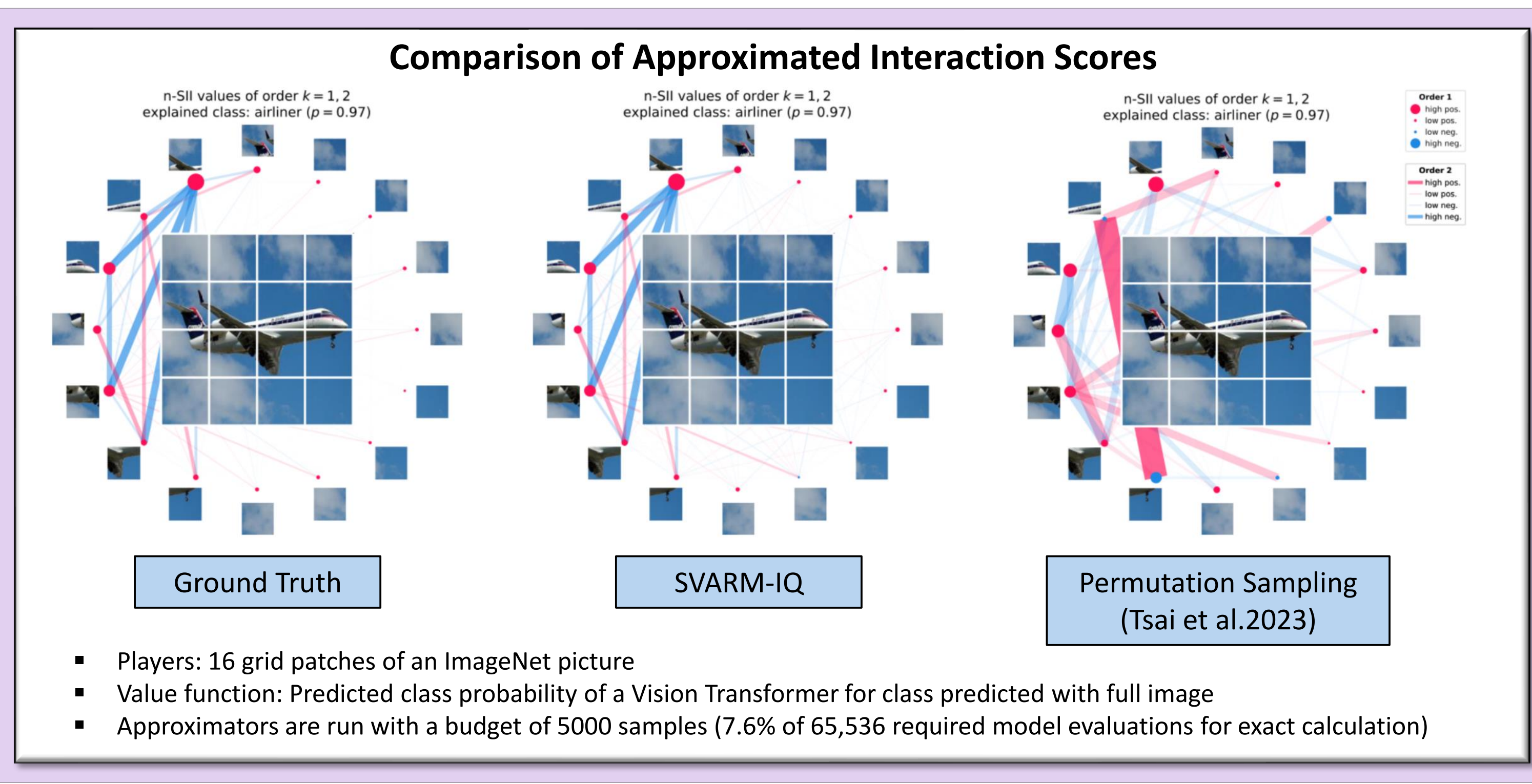
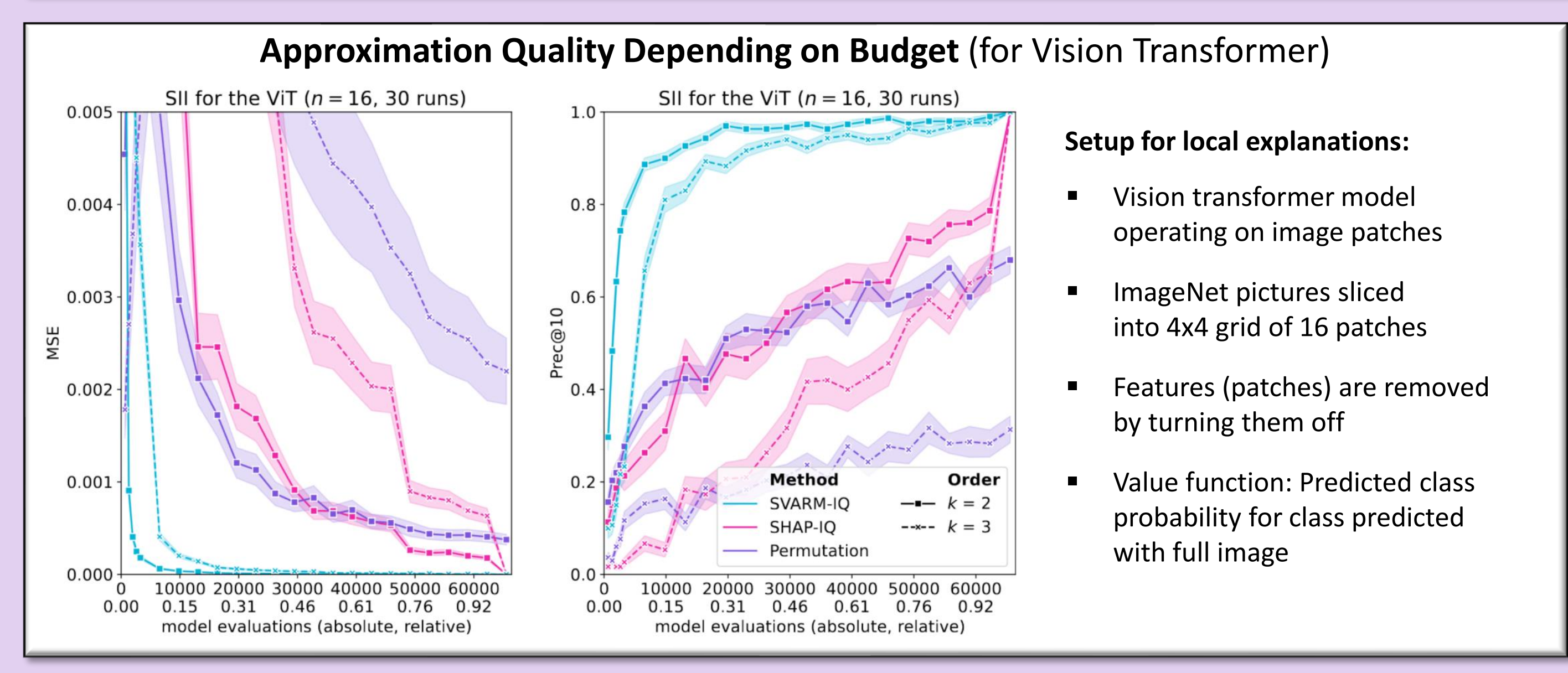
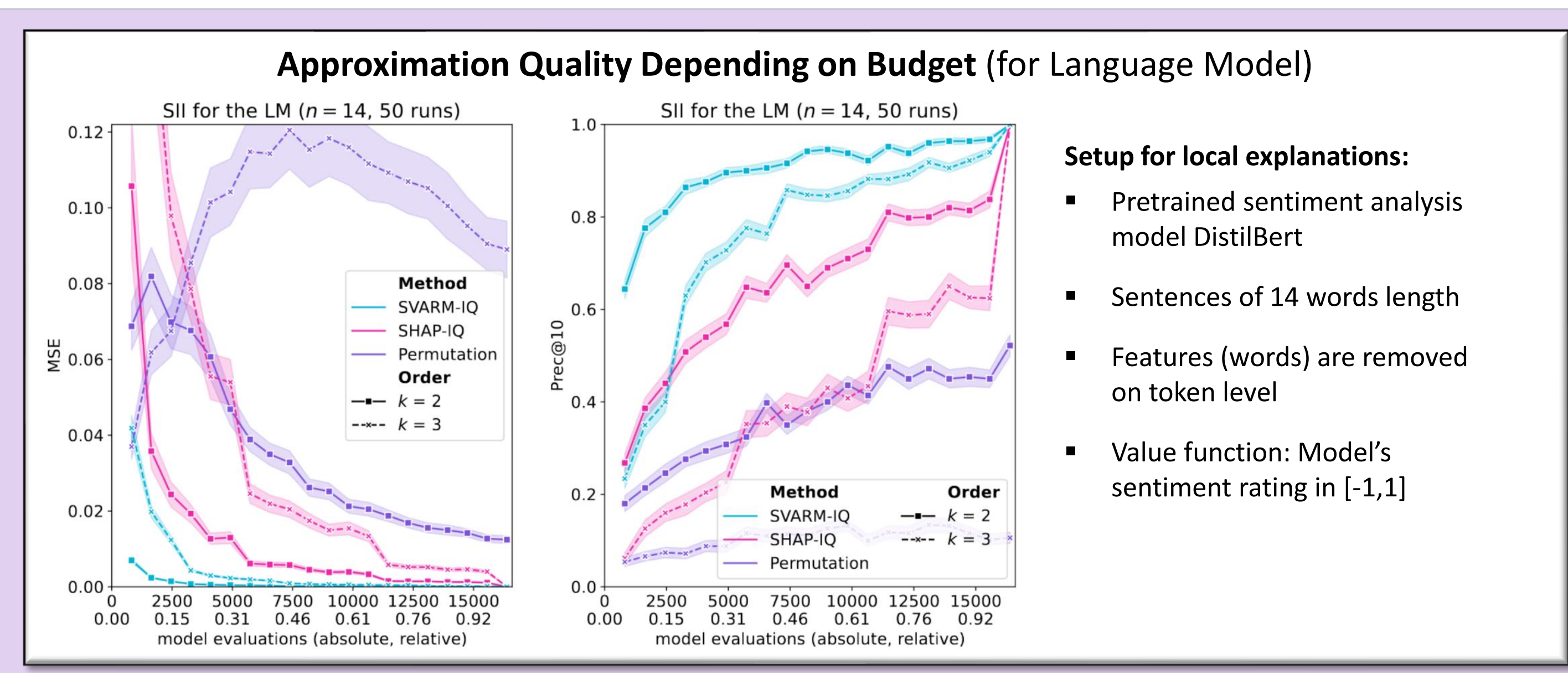
- Sampling Coalitions and Updating all Estimates**
- Calculate all border strata exactly with coalitions of size $0, \dots, b, n-b, \dots, n$
 - Perform warmup on inner strata $b+1, \dots, n-b-1$ to initialize all estimates
 - Repeat with remaining budget \tilde{B} :
 - Draw coalition size $s \in \{b+1, \dots, n-b-1\} \sim \tilde{P}(s)$
 - Draw coalition A of size s uniformly at random and evaluate $\nu(A)$
 - For all K : Update $I_{K,\ell}^W$ with $W = A \cap K$ and $\ell = |A| - |W|$
-

Theorem 4.2 & Corollary 4.3: Variance and MSE

For any K of order k , the variance and MSE of the estimate \hat{I}_K returned by SVARM-IQ is bounded by

$$\mathbb{E} [(\hat{I}_K - I_K)^2] = \mathbb{V}[\hat{I}_K] \leq \frac{\gamma_k}{\tilde{B}} \sum_{W \subseteq K} \sum_{\ell=0}^{n-k} \binom{n-k}{\ell}^2 \lambda_{k,\ell}^2 \sigma_{K,\ell,W}^2$$

- With stratum variances $\sigma_{K,\ell,W}^2 = \mathbb{V}[\nu(A \cup W)]$ for $A \subseteq \mathcal{N} \setminus K$ with $|A| = \ell$ drawn u.a.r.
- $\gamma_k := 2(n-1)^2$ for $k=2$ and $\gamma_k = n^{k-1}(n-k+1)^2$ for $k \geq 3$



References

Fumagalli, F., Muschalik, M., Kolpaczki, P., Hüllermeier, E., and Hammer, B. (2023). SHAP-IQ: Unified Approximation of any-Order Shapley Interactions. In *Proceedings of Advances in Neural Processing Systems*.

Fujimoto, K., Kojadinovic, I., and Marichal, J. (2006). Axiomatic characterizations of probabilistic and cardinal-probabilistic interaction indices. *Games and Economic Behavior*, 55(1):72-99.

Grabisch, M. and Roubens, M. (1999). An axiomatic approach to the concept of interaction among players in cooperative games. *International Journal of Game Theory*, 28(4):547-565.

Kolpaczki, P., Bengs, V., Muschalik, M., Hüllermeier, E. (2024). Approximating the Shapley Value without Marginal Contributions. In *Proceedings of AAAI Conference on Artificial Intelligence*.

Shapley, L.S. (1953). A Value for n -Person Games. In *Contributions to the Theory of Games (AM-28)*, Volume II, pages 307-318. Princeton University Press.

Sundararajan, M., Dharmdhere, K., and Agrawal, A. (2020). The Shapley Taylor Interaction Index. In *Proceedings of the 37th International Conference on Machine Learning*.

Tsai, C., Yeh, C., and Ravikumar, P. (2023). Faith-Shap: The Faithful Shapley Interaction Index. *Journal of Machine Learning Research*, 24(94):1-42.