

How much data is sufficient to learn high-performing algorithms?

Generalization guarantees for data-driven algorithm design

Nina Balcan, Dan DeBlasio, Travis Dick, Carl Kingsford, Tuomas Sandholm, **Ellen Vitercik**

Data-driven algorithm design

Algorithms often have **many tunable parameters**
Significant impact on runtime, solution quality, ...

Hand-tuning is **time-consuming, tedious**, and **error prone**

Goal: Automate algorithm configuration via machine learning

Input: **Training set** of typical problem instances from application at hand

Sampled from unknown, **application-specific distribution**

Output: Configuration with strong **average empirical performance** on training set

Runtime, solution quality, etc.

Parameter setting should—**ideally**—be good on future inputs

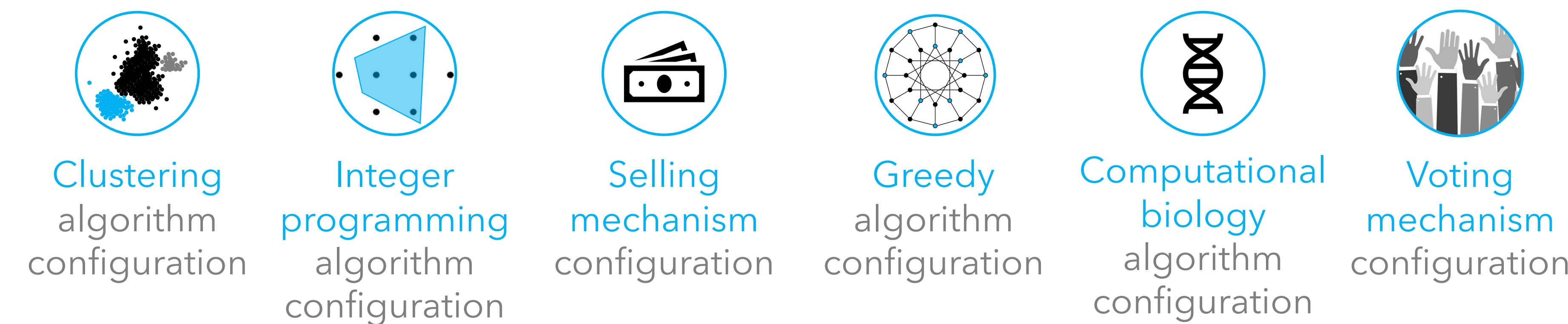
Summary of contributions

Broadly applicable theory for deriving **generalization bounds**:

$$\left[\begin{array}{c} \text{Algorithm's average performance} \\ \text{on training set} \end{array} \right] - \left[\begin{array}{c} \text{Algorithm's expected performance} \\ \text{on unknown distribution} \end{array} \right] \leq ?$$

Prior research proved generalization bounds case-by-case

Gupta, Roughgarden, ITCS'16; Balcan, Nagarajan, V, White, COLT'17; Balcan, Dick, Sandholm, V, ICML'18; Balcan, Dick, White, NeurIPS'18; Balcan, Dick, Lang, ICLR'20; ...



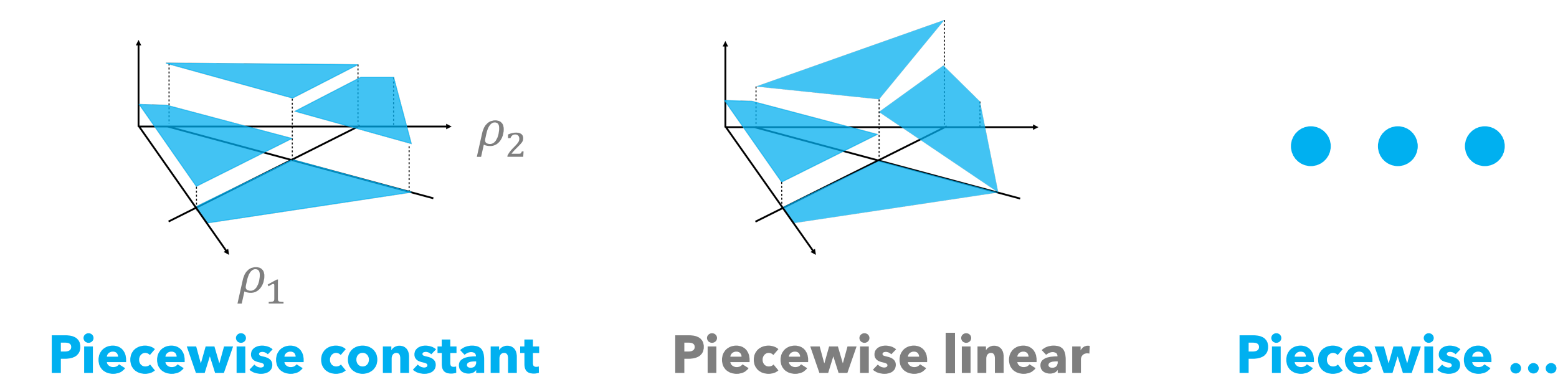
Recover bounds

Prove novel bounds

We uncover **overarching structure** linking these seemingly disparate domains

Guarantees apply to any parameterized algorithm where:

Performance is a **piecewise-structured** function of parameters

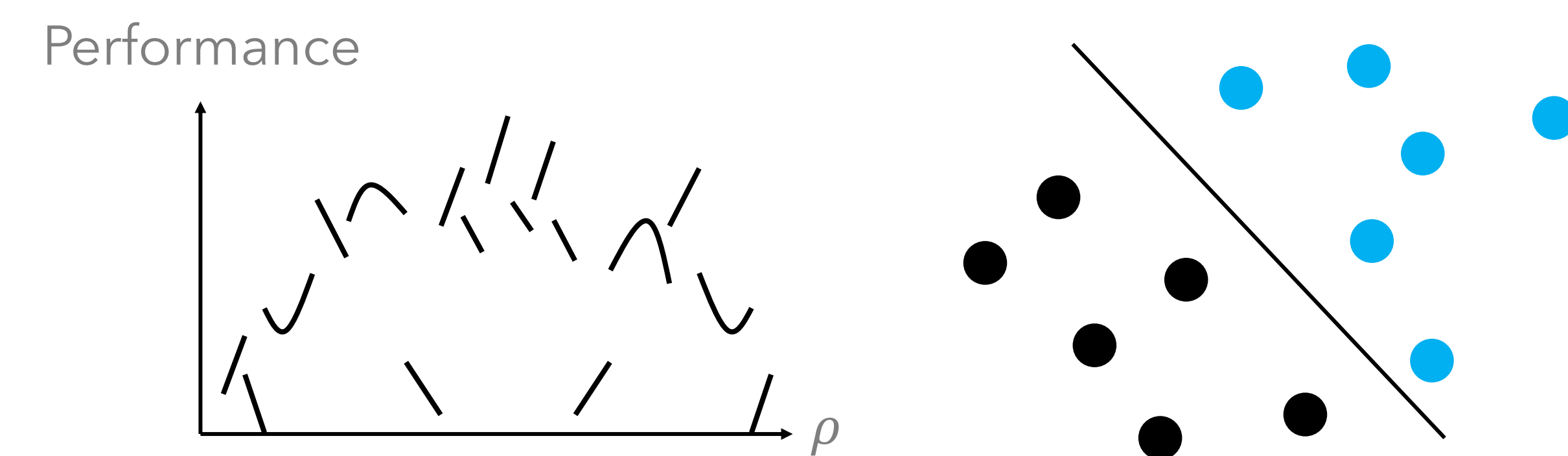


Additional references

- Book chapter by Balcan [Cambridge University Press '20]
- Online algorithm configuration:
Exploited that the dual functions are piecewise Lipschitz to provide regret bounds
[Balcan, Dick, V, FOCS'18; Balcan, Dick, Pegden, UAI'20; Balcan, Dick, Sharma, AISTATS'20]

Primary challenge

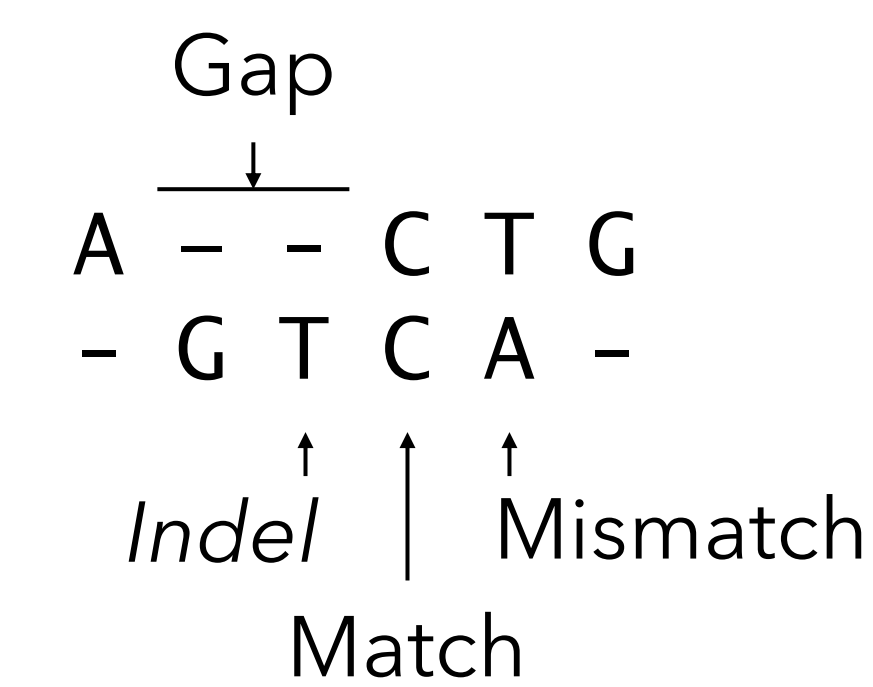
Performance is a **volatile** function of parameters
Complex connection between parameters and performance



Meanwhile, for well-understood functions in machine learning theory:

Simple connection between function parameters and value

Running example: Sequence alignment



Standard algorithm with parameters $\rho_1, \rho_2, \rho_3 \geq 0$:

Return alignment maximizing:
(# matches) $- \rho_1 \cdot$ (# mismatches) $- \rho_2 \cdot$ (# indels) $- \rho_3 \cdot$ (# gaps)

"There is **considerable disagreement** among molecular biologists about the **correct choice** [of ρ]" [Gusfield et al. '94]

Model and problem formulation

\mathbb{R}^d : Set of all parameters

\mathcal{X} : Set of all inputs (e.g., sequence pairs)

$u_\rho(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
Runtime, solution quality, ...

Assume $u_\rho(x) \in [-1, 1]$

Standard assumption: Unknown distribution \mathcal{D} over inputs
Models specific application domain at hand

Generalization bound: Given samples $x_1, \dots, x_N \sim \mathcal{D}$, for any ρ ,

$$\left| \frac{1}{N} \sum_{i=1}^N u_\rho(x_i) - \mathbb{E}_{x \sim \mathcal{D}} [u_\rho(x)] \right| \leq ?$$

Empirical average utility **Expected utility**

Main result

$\mathcal{U} = \{u_\rho: \mathcal{X} \rightarrow \mathbb{R} \mid \rho \in \mathbb{R}^d\}$ **"Primal" function class**

Typically, prove generalization guarantees by bounding the **complexity** of \mathcal{U}
VC dimension, Rademacher complexity, ...

Challenge: \mathcal{U} is gnarly. E.g., in sequence alignment:

- Each domain element is a pair of sequences
- Unclear how to plot/visualize functions u_ρ
- No obvious notions of Lipschitzness or smoothness to rely on

This is where **dual functions** come in handy!

$u_x^*(\rho)$ = utility as function of parameters

$u_x^*(\rho) = u_\rho(x)$

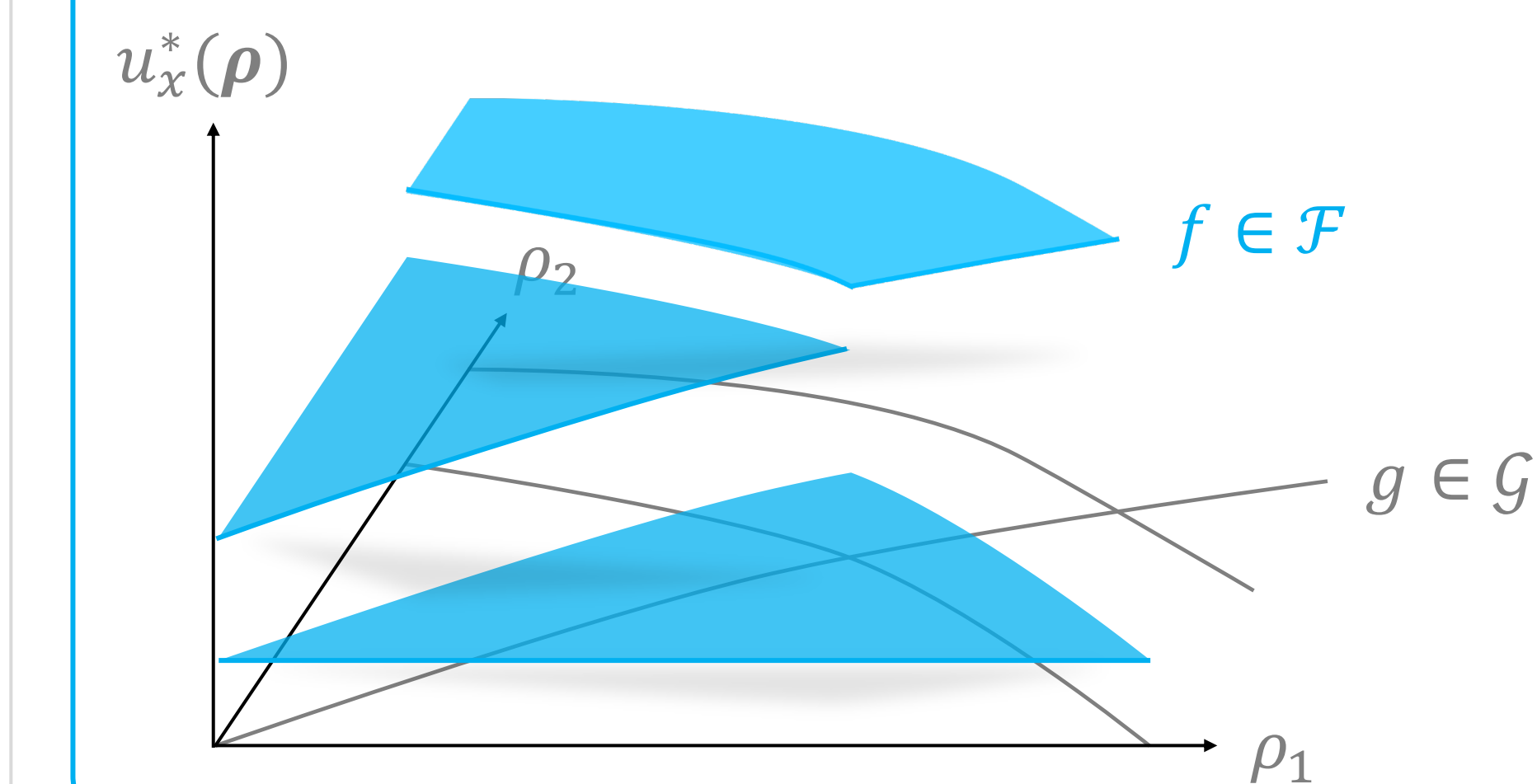
$\mathcal{U}^* = \{u_x^*: \mathbb{R}^d \rightarrow \mathbb{R} \mid x \in \mathcal{X}\}$ **"Dual" function class**

Across algorithm configuration, ubiquitously, the duals are **piecewise-structured**

Theorem

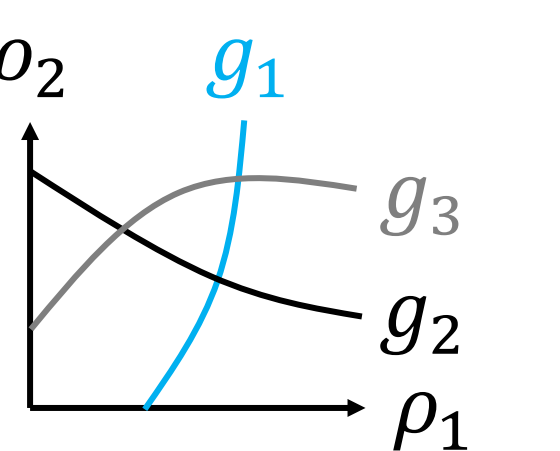
With high probability, for all ρ :

$$|\text{Average utility on training set} - \text{expected utility}| = \tilde{O} \left(\frac{\sqrt{\frac{\text{\# boundary functions} \cdot (\mathcal{P} - \text{dim}(\mathcal{F}^*) + \text{VC-dim}(\mathcal{G}^*) \ln k)}{N}}}{N} \right)$$



Lemma: Given k boundaries, how many sign patterns do they make?

$$\left| \left\{ \begin{pmatrix} g_1(\rho) \\ \vdots \\ g_k(\rho) \end{pmatrix} : \rho \in \mathbb{R}^d \right\} \right| \leq (ek)^{\text{VCdim}(\mathcal{G}^*)}$$



Proof idea: Transition to dual and apply Sauer's lemma: for any ρ_1, \dots, ρ_k

$$\left| \left\{ \begin{pmatrix} g(\rho_1) \\ \vdots \\ g(\rho_k) \end{pmatrix} : g \in \mathcal{G} \right\} \right| \leq (ek)^{\text{VCdim}(\mathcal{G})}$$

Example application: Sequence alignment

With high probability, for any $\rho \in \mathbb{R}^3$,

$$|\text{avg utility on training set} - \text{expected utility}| = \tilde{O} \left(\sqrt{\frac{\ln(\text{seq. length})}{N}} \right)$$

Distance between algorithm's output and ground-truth alignment