

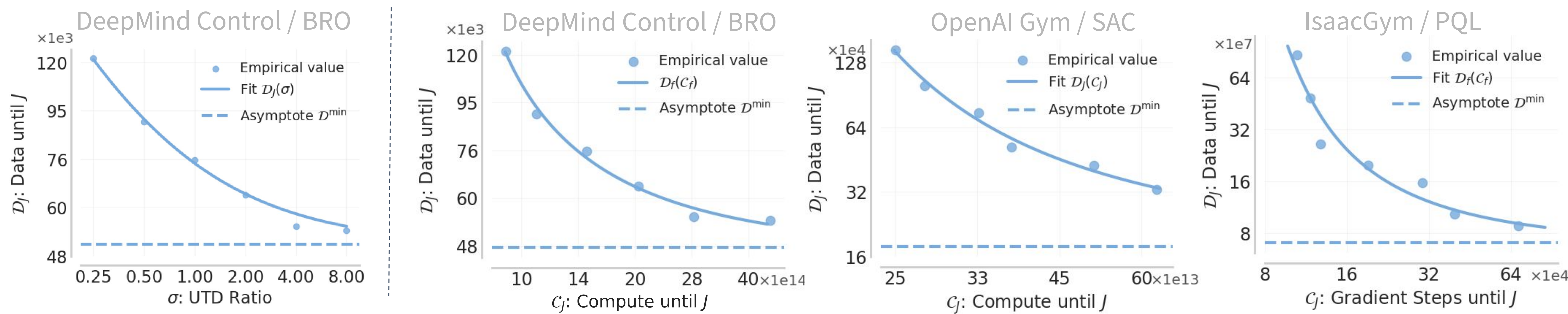
Motivation and Problem Statement

- Q1:** RL is sensitive to hyperparameters. How to set batch size and learning rate for large-scale runs?
- Q2:** To achieve performance level J , RL requires a combination of *data* \mathcal{D} and *compute* \mathcal{C} . What is their tradeoff?
- Q3:** I have a requirement on the total *budget*, a combination $\mathcal{F} = \mathcal{C} + \delta \cdot \mathcal{D}$ of data and compute. How should I configure the algorithm to maximize performance given this budget?

Q2, Q3: Configuring Algorithm

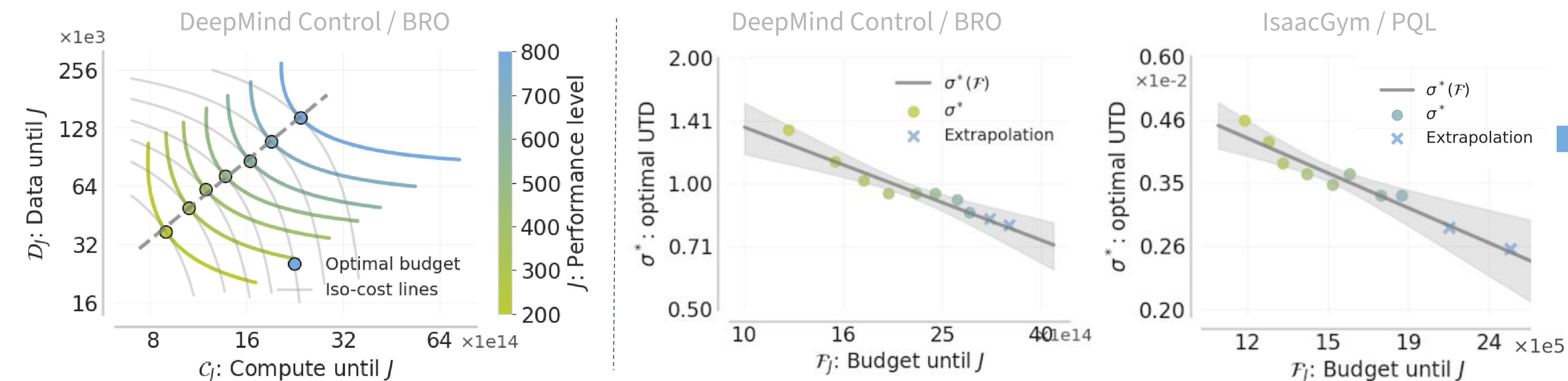
Q2 With the right hyperparameters, data and compute are predictable functions of UTD! Data follows a power law across environments/algorithms.

Tradeoff along the data-compute Pareto frontier parameterized by UTD.



Q3 Each J corresponds to a $\mathcal{D}-\mathcal{C}$ Pareto frontier. Each frontier has a unique budget-optimal point, corresponding to an optimal UTD* and budget \mathcal{F} .

Optimal UTD follows a power law wrt budget!

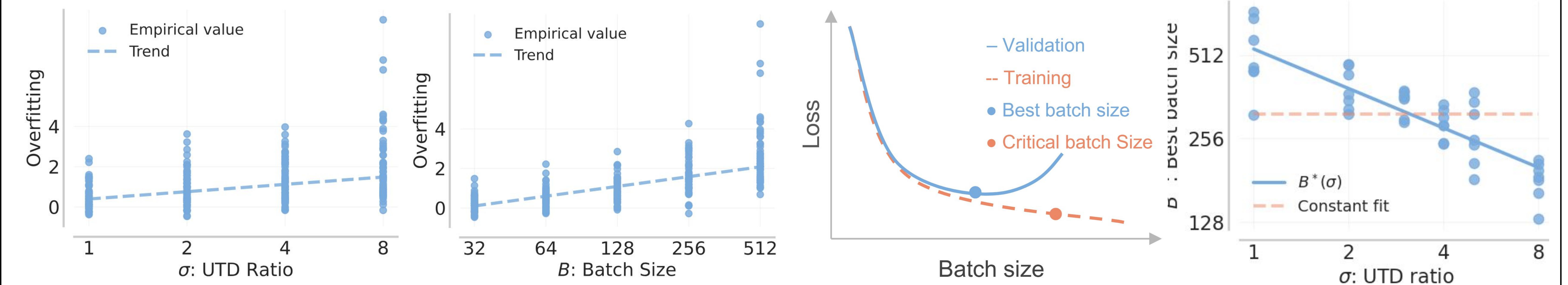


Q1: Setting Hyperparameters

Batch size $\text{Overfitting} = (\text{TD error on newest data}) - (\text{TD error on uniformly sampled data})$.

High *updates-to-data ratio* (UTD) and batch size both increase overfitting, since they train on the average data point more. **Decrease batch size with higher UTD to counteract overfitting.**

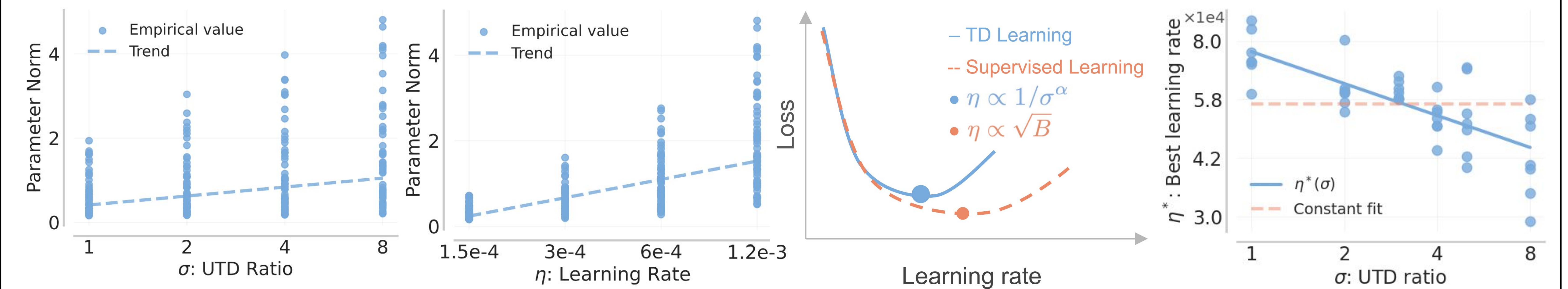
$$B^*(\text{UTD}) = \left(\frac{\beta_B}{\text{UTD}} \right)^{\alpha_B}$$



Learning rate High UTD ratios and high learning rates both lead to *plasticity loss*.

Decrease learning rate to counteract high parameter norm.

$$\eta^*(\text{UTD}) \approx \left(\frac{\beta_{\eta}}{\text{UTD}} \right)^{\alpha_{\eta}}$$



Summary and Takeaways

- Batch size, learning rate predictable from UTD; outperforms constant.
- UTD controls the tradeoff between data, compute.
- Optimal UTD follows a power law wrt budget.

