



Motivation

Problem: Target-domain labeled data are **scarce**; source-domain data are abundant but **shifted**.

We observe the following two categories of data:

- ▶ Target-domain data $\mathcal{D}_T = \{(X_i^{(0)}, Y_i^{(0)})\}_{i=1}^{n_0} \sim P^{(0)}$.
- ▶ K source-domain data $\mathcal{D}_S^{(k)} = \{(X_i^{(k)}, Y_i^{(k)})\}_{i=1}^{n_k} \sim P^{(k)}$ for $k = 1, \dots, K$.

Two standard assumptions in the transfer learning literature:

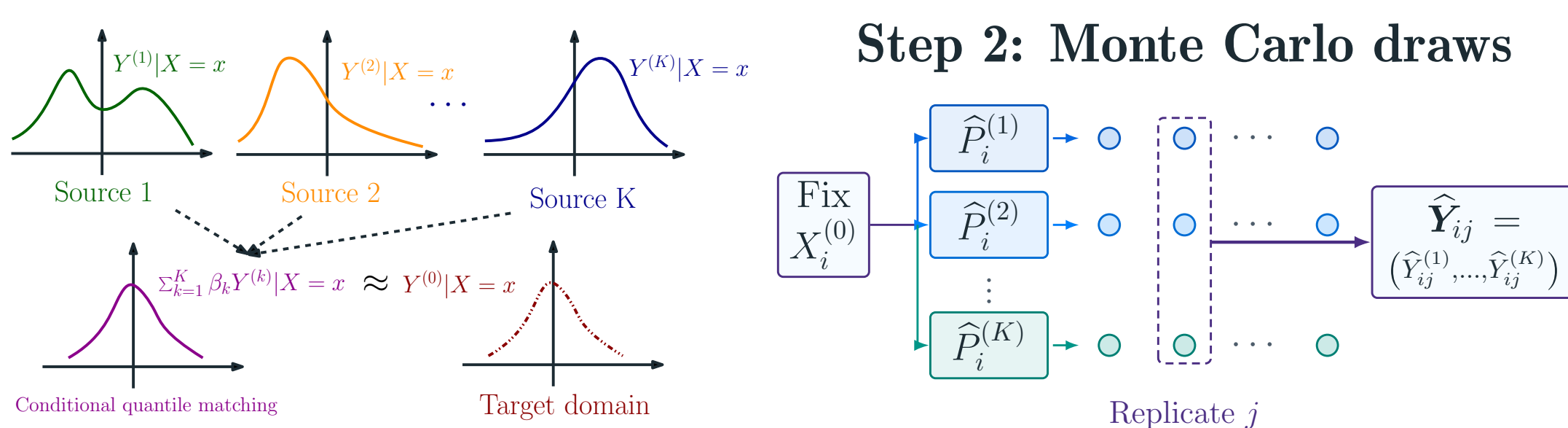
- ▶ *Covariate Shift:* $P^{(k)}(y|x) = P^{(0)}(y|x)$ but $P^{(k)}(x) \neq P^{(0)}(x)$;
- ▶ *Label Shift:* $P^{(k)}(x|y) = P^{(0)}(x|y)$ but $P^{(k)}(y) \neq P^{(0)}(y)$.

Synthetic Data Augmentation: Generate high-quality labeled data for the target domain when, for all $k = 1, \dots, K$,

$$P^{(0)}(x) \neq P^{(k)}(x) \quad \text{and} \quad P^{(0)}(y|x) \neq P^{(k)}(y|x).$$

TLCQM Pipeline

Main Idea: Learn source conditional generators, generate synthetic responses at target covariates, then calibrate them by conditional quantile matching.



Step 1: Learn a generator $\hat{P}^{(k)}(y|x)$ for each source domain $\mathcal{D}_S^{(k)}$.

Step 2: For target covariates $X_i^{(0)}$, draw synthetic response vectors $\hat{Y}_{ij} = (\hat{Y}_{ij}^{(1)}, \dots, \hat{Y}_{ij}^{(K)})$, $j = 1, \dots, M$ with $\hat{Y}_{ij}^{(k)} \sim \hat{P}^{(k)}(y|X_i^{(0)})$.

Step 3: Compute the conditional quantile matching estimator

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^{K+1}} \sum_{i=1}^{n_0} \sum_{j=1}^M \left[Y_{ij}^{(0)} - (\beta^T \hat{V})_{(ij)} \right]^2, \quad \hat{V}_{ij} = (1, \hat{Y}_{ij})^T,$$

where $Y_{ij}^{(0)}$ and $(\beta^T \hat{V})_{(ij)}$ are the corresponding order statistics.

Step 4 (optional): Estimate $w_k(x) = \frac{dP^{(0)}(x)}{dP^{(k)}(x)}$ for covariate shift.

Step 5: Obtain the augmented target data $\mathcal{D}_T \cup \mathcal{D}_A$, where $\mathcal{D}_A = \bigcup_{k=1}^K \left\{ (X_i^{(k)}, \beta^T \hat{V}_i^{(k)}) \right\}_{i=1}^{n_k}$ with $\hat{V}_i^{(k)} = (1, \hat{Y}_i^{(1,k)}, \dots, \hat{Y}_i^{(K,k)})$ and $\hat{Y}_i^{(j,k)}$ as any predicted value for $Y^{(j)}$ from $\hat{P}^{(j)}(y|X_i^{(k)})$.

Theory: Excess Risk Decomposition

With rich source data, TLCQM moves the leading ERM complexity term from n_0 to $N = \sum_{k=0}^K n_k$, while making other error dependencies explicit.

- ▶ $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a loss function, and the population risk of a prediction function $f \in \mathcal{F}$ is $R(f) := \mathbb{E}_{P^{(0)}} [\ell(Y^{(0)}, f(X^{(0)}))]$.
- ▶ $\text{Rad}_n(\mathcal{F}) = \mathbb{E} \left\{ \frac{1}{n} \cdot \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i \cdot f(X_i) \right| \right] \right\}$ is the Rademacher complexity of \mathcal{F} .
- ▶ \mathcal{B} is the solution set of population-level conditional quantile matching.

Target-only Excess Risk: Let $\hat{f}^{(0)}$ be the empirical risk minimizer (ERM) under target-only data:

$$R(\hat{f}^{(0)}) - R(f^{(0)}) \lesssim \text{Rad}_{n_0}(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{n_0}} \quad \text{w.p.} \geq 1 - \delta.$$

Excess Risk Under TLCQM: Define the transfer learning prediction function under our TLCQM framework as $\hat{f}^{(0,t)} = \arg \min_{f \in \mathcal{F}} \frac{1}{N} \left\{ \sum_{i=1}^{n_0} \ell(Y_i^{(0)}, f(X_i^{(0)})) + \sum_{k=1}^K \sum_{i=1}^{n_k} \hat{w}_k(X_i^{(k)}) \cdot \ell(\hat{\beta}^T \hat{V}_i^{(k)}, f(X_i^{(k)})) \right\}$, leading to the bound

$$R(\hat{f}^{(0,t)}) - R(f^{(0)}) \lesssim \underbrace{\text{Rad}_N(\mathcal{F})}_{\text{standard generalization error}} + \underbrace{\frac{1}{N} \sum_{k=1}^K \|\hat{w}_k - w_k\|_1}_{\text{importance weight error}} + \underbrace{\inf_{\beta \in \mathcal{B}} \|\hat{\beta} - \beta_*\|_1}_{\text{quantile matching error}} + \underbrace{\sum_{k=1}^K \|\hat{g}^{(k)} - g^{(k)}\|_\infty}_{\text{distributional learning error}} + \underbrace{\inf_{\beta \in \mathcal{B}} \left(\int_0^1 [Q_{Y^{(0)}}(\alpha) - Q_{\beta^T V}(\alpha)]^2 d\alpha \right)^{1/2}}_{\text{transfer bias}}.$$

- ▶ Our result is **learner-agnostic**: transfer quality is separated from the downstream function class through the five explicit error terms.

Theory: Quantile Matching Error

$$\inf_{\beta \in \mathcal{B}} \|\hat{\beta} - \beta_*\|_1 = O \left(\sqrt{\frac{K}{n_0} \sum_{k=1}^K \|\hat{g}^{(k)} - g^{(k)}\|_\infty} \right) + O \left(\frac{1}{\sqrt{M}} \right) + O_P \left(\sqrt{\frac{K \log \log n_0}{n_0}} + \sqrt{K} \left[\frac{\log \log n_0}{n_0} \inf_{\beta \in \mathcal{B}} \int_0^1 \{Q_{Y^{(0)}}(\alpha) - Q_{\beta^T V}(\alpha)\}^2 d\alpha \right]^{1/4} \right).$$

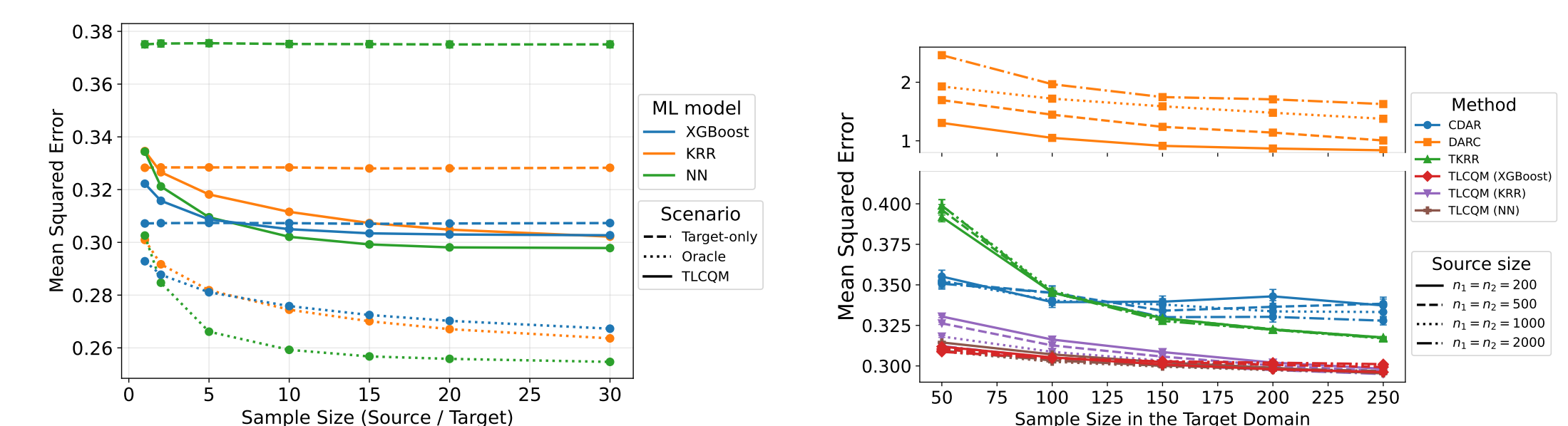
- ▶ The rate separates generator error, Monte Carlo error, target-size stochastic error, and a product of **transfer bias** and $O_P(n_0^{-1/4})$.
- ▶ **Transfer bias** vanishes when the target conditional response distribution lies in the convex hull of source conditional response distributions.

Small transfer bias improves both approximation and the convergence rate of the quantile matching objective.

Simulations: General Target Shift

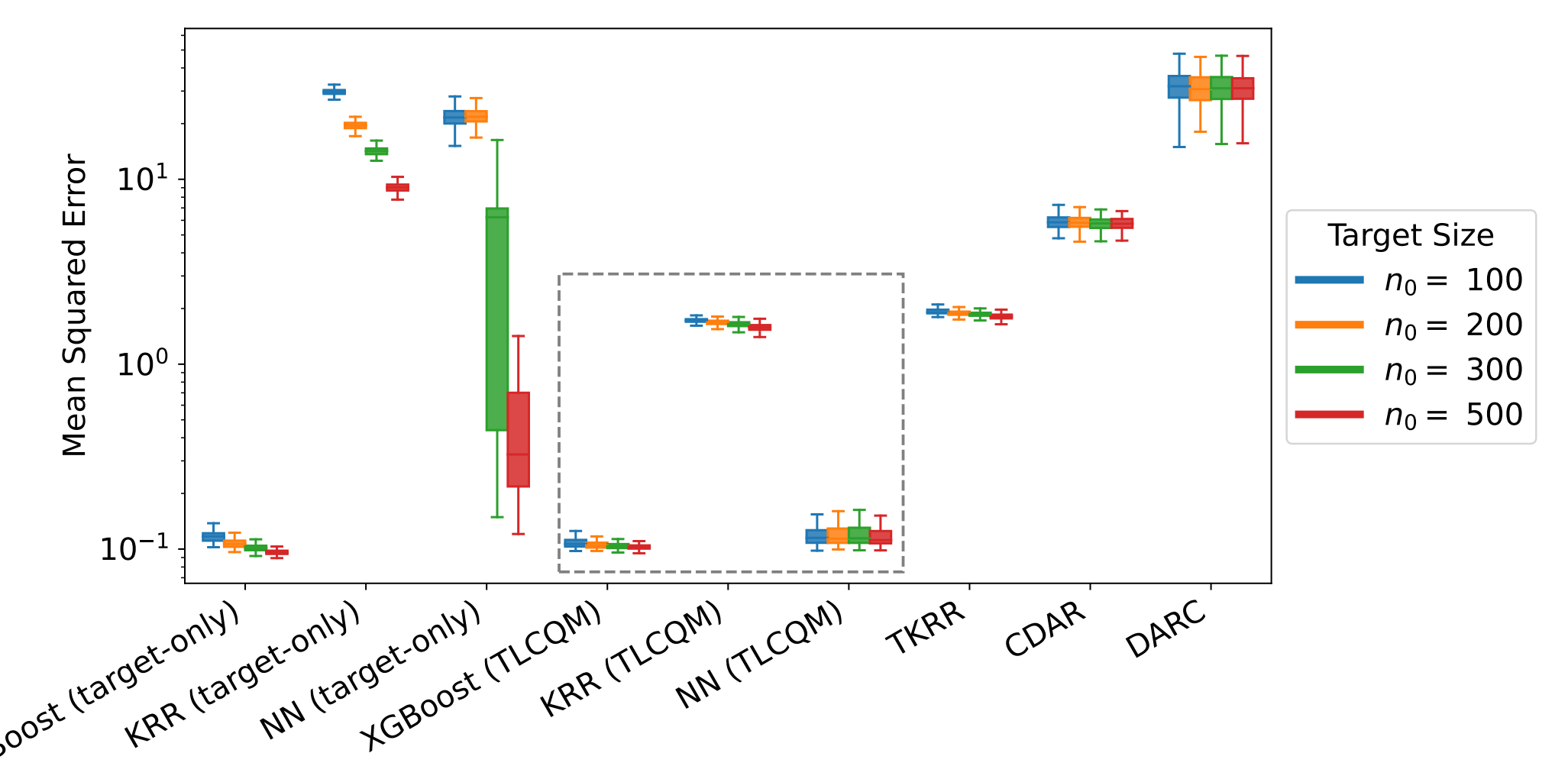
We model each source-domain conditional distribution as $g^{(k)}(x, \cdot) \# P_\eta = P^{(k)}(\cdot|x)$, and estimate $g^{(k)}$ via engression: $\hat{g}^{(k)} \in \arg \min_{g \in \mathcal{F}} \frac{1}{n_k} \sum_{i=1}^{n_k} \left[\frac{1}{m} \sum_{j=1}^m |Y_i^{(k)} - g(X_i^{(k)}, \eta_{ij})| - \frac{1}{2m(m-1)} \sum_{j=1}^m \sum_{j'=1}^m |g(X_i^{(k)}, \eta_{ij}) - g(X_i^{(k)}, \eta_{ij'})| \right]$, where P_η is a pre-specified noise distribution.

- ▶ Source domains: $Y^{(1)} = \sin(3\theta^T X^{(1)}) + 1 + \epsilon$ and $Y^{(2)} = 2 \cos(3\theta^T X^{(2)}) + 1 + \epsilon$, where $\theta = (1, \frac{1}{2}, \dots, \frac{1}{6})^T \in \mathbb{R}^6$ and $X^{(1)}, X^{(2)} \sim \mathcal{N}(\mathbf{1}_6, \mathbf{I}_6)$, $\epsilon \sim \mathcal{N}(0, 0.25)$ with $\mathbf{1}_6 = (1, \dots, 1)^T \in \mathbb{R}^6$ and $\mathbf{I}_6 \in \mathbb{R}^{6 \times 6}$ being the identity matrix.
- ▶ Target domain: $Y^{(0)} = \frac{1}{3} \sin(3\theta^T X^{(0)}) - 3 + \epsilon$ with $X^{(0)} \sim \mathcal{N}(\mathbf{0}_6, 0.25 \cdot \mathbf{I}_6)$ and $\epsilon \sim \mathcal{N}(0, 0.25)$.

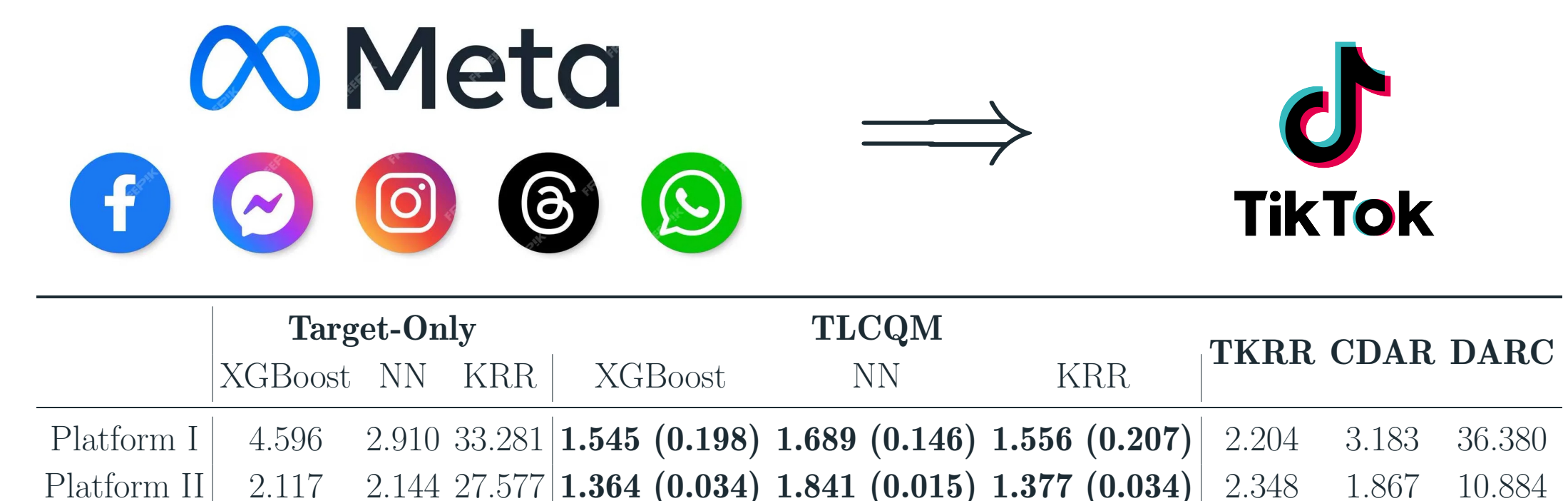


Real Data Applications

Public Apartment Data: Predict log (apartment rental price) for the target state FL using data from source states IL, OH, and WA.



Monthly Active User Prediction at Meta: Predict the country-level monthly active user statistics of a target app using data from Meta's family of apps.



- ▶ Flexible learners, especially KRR and NN, benefit most from our TLCQM framework when the target size n_0 is small.