

**ICML**  
International Conference  
On Machine Learning

# Heterogeneous Data Game: Characterizing the Model Competition Across Multiple Data Sources

**Renzhe Xu, Kang Wang, Bo Li**

Shanghai University of Finance and Economics, Tianjin University, Tsinghua University

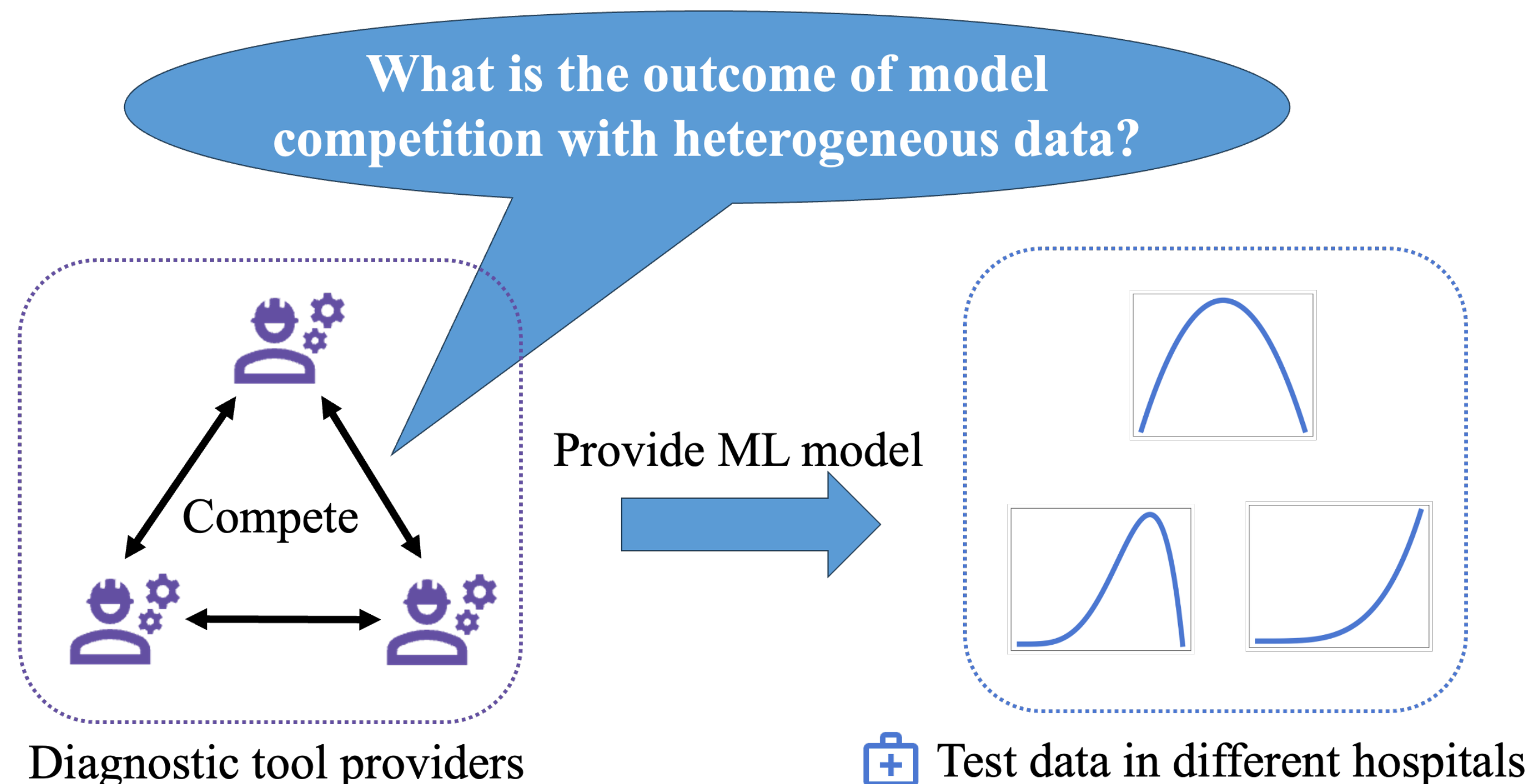
# Background

- **Model competition with heterogeneous data**
  - **Data heterogeneity:** Machine learning (ML) models face test data from diverse data sources with varying distributions
  - **ML model competition:** Multiple ML model providers compete for market share across these data sources

# Background

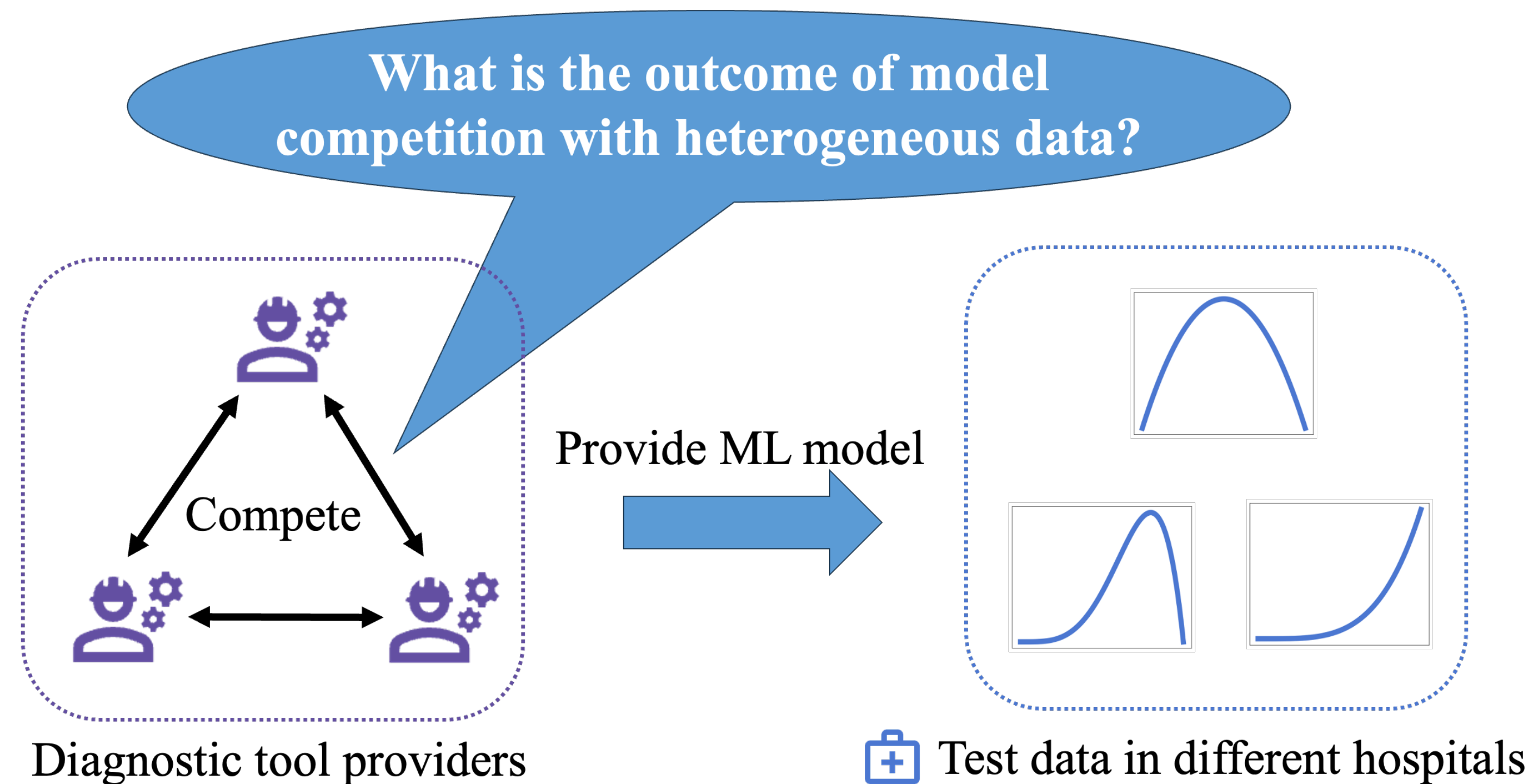
- **Example in health care**

- Different hospitals have different patient populations
- Multiple diagnostic tool providers compete across hospitals



# Background

- **Our goal**
  - What equilibrium arises in such competition?
  - What factors shape different kinds of equilibrium?



# Problem setup

- **Basic parameters**
  - $N$  model providers and  $K$  data sources
  - Each data source  $k \in [K]$  has an importance weight  $w_k$  and a data distribution  $P_k(x, y)$

# Problem setup

- **Basic parameters**

- $N$  model providers and  $K$  data sources
- Each data source  $k \in [K]$  has an importance weight  $w_k$  and a data distribution  $P_k(x, y)$

**Assumption:**

$P_k(y|x)$  is a linear model

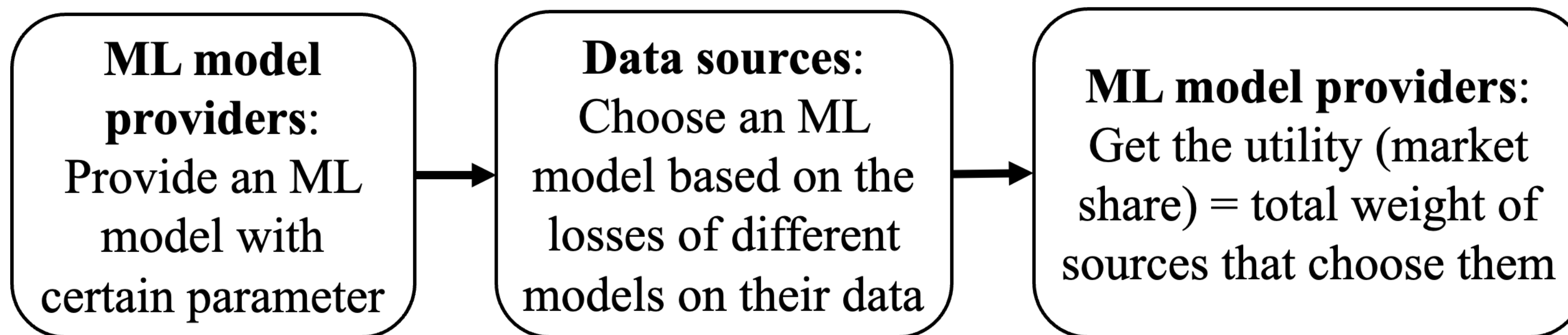
# Problem setup

- **Basic parameters**

- $N$  model providers and  $K$  data sources
- Each data source  $k \in [K]$  has an importance weight  $w_k$  and a data distribution  $P_k(x, y)$

**Assumption:**  
 $P_k(y|x)$  is a linear model

- **Heterogeneous data game**



# Our results

- **Three patterns of pure Nash equilibrium (PNE)**
  - **Non-existence of PNE**
  - **Homogeneous PNE**
    - All providers select the same ML model
  - **Heterogeneous PNE**
    - Providers offer different ML models and specialize in different data sources



# Our results

|  | <b>Proximity choice model</b><br>(Data sources choose the model with the lowest loss)  | <b>Probability choice model</b><br>(Data sources choose models via a loss-based logit function)   |
|--|--|---|
| <b>Monopoly</b><br><b>(<math>N = 1</math>)</b> | Model provider minimize the $w_k$ -weighted loss over all data sources   |   |
| <b>Duopoly</b><br><b>(<math>N = 2</math>)</b>  | <ul style="list-style-type: none"> <li>• Equivalent condition for PNE existence</li> <li>• PNE must be <b>heterogeneous</b>, if it exists</li> </ul> | <ul style="list-style-type: none"> <li>• Equivalent condition for PNE existence</li> <li>• PNE must be <b>homogeneous</b>, if it exists</li> </ul>  |
| <b><math>N &gt; 2</math></b>                   | <ul style="list-style-type: none"> <li>• Sufficient condition for PNE existence</li> <li>• PNE must be <b>heterogeneous</b>, if it exists</li> </ul> | <ul style="list-style-type: none"> <li>• Equivalent condition for <b>homogeneous</b> PNE existence</li> <li>• Sufficient condition for <b>heterogeneous</b> PNE existence</li> <li>• Example when both types of PNE exist simultaneously</li> </ul> |

# Conclusions

- **Our contribution**

- Model the competition among model providers in a heterogeneous data environment
- Study how market factors affect the resulting equilibrium
  - E.g., the number of model providers, data source choice model

- **Potential policy implications**

- Market designers could adjust market factors to achieve a more desirable equilibrium

**Thanks for listening!**

Renzhe Xu, Shanghai University of Finance and Economics

Paper: <https://arxiv.org/abs/2505.07688>

Email: [xurenzhe@sufe.edu.cn](mailto:xurenzhe@sufe.edu.cn)