# Bridging Structural Causal Inference and Machine Learning: The S-DIDML Estimator for Heterogeneous Treatment Effects

**Yile Yu[1]\*, Anzhi Xu[2]**

1.School of Education, Zhejiang University of Technology, Hangzhou, 310014, China

2.School of Accounting, Yunnan University of Finance and Economics, Kunming, 650221, China

*\*Corresponding author: Yile Yu, 302023572139@zjut.edu.cn*

**Abstract:** In response to the increasing complexity of policy environments and the proliferation of high-dimensional data, this paper introduces the S-DIDML estimator—a framework grounded in structure and semiparametrically flexible for causal inference. By embedding Difference-in-Differences (DID) logic within a Double Machine Learning (DML) architecture, the S-DIDML approach combines the strengths of temporal identification, machine learning-based nuisance adjustment, and orthogonalized estimation. We begin by identifying critical limitations in existing methods, including the lack of structural interpretability in ML models, instability of classical DID under high-dimensional confounding, and the temporal rigidity of standard DML frameworks. Building on recent advances in staggered adoption designs and Neyman orthogonalization, S-DIDML offers a five-step estimation pipeline that enables robust estimation of heterogeneous treatment effects (HTEs) while maintaining interpretability and scalability. Demonstrative applications are discussed across labor economics, education, taxation, and environmental policy. The proposed framework contributes to the methodological frontier by offering a blueprint for policy-relevant, structurally interpretable, and statistically valid causal analysis in complex data settings.
**Keywords:** S-DIDML; Methodology; Causal Inference; Difference-in-Differences; Double Machine Learning; Semiparametric Methods

## 1.Introduction

In the social sciences, the pursuit of causal understanding—rather than mere correlation—has long defined the discipline's scientific ambition. Researchers in various fields endeavor to isolate the effect of a particular treatment or intervention from the myriad of other confounding influences. This is true when assessing the effect of a trade policy on employment, the impact of educational reform on student performance, or the consequences of environmental regulations on firm productivity. At the core of these efforts lies causal inference, which seeks to address the fundamental question: It is imperative to consider what the outcome would have been in the absence of the treatment. Economics, political science, education, and public policy have increasingly embraced quasi-experimental designs to approximate this elusive counterfactual.[1] Among them, the Difference-in-Differences (DID) framework has become a canonical tool.[2] It leverages temporal and group-based variation to identify causal effects under a critical "parallel trends" assumption.[3] The appeal of DID lies in its structural interpretability, transparency, and policy relevance, making it a mainstay in top journals and institutional evaluations alike. However, the
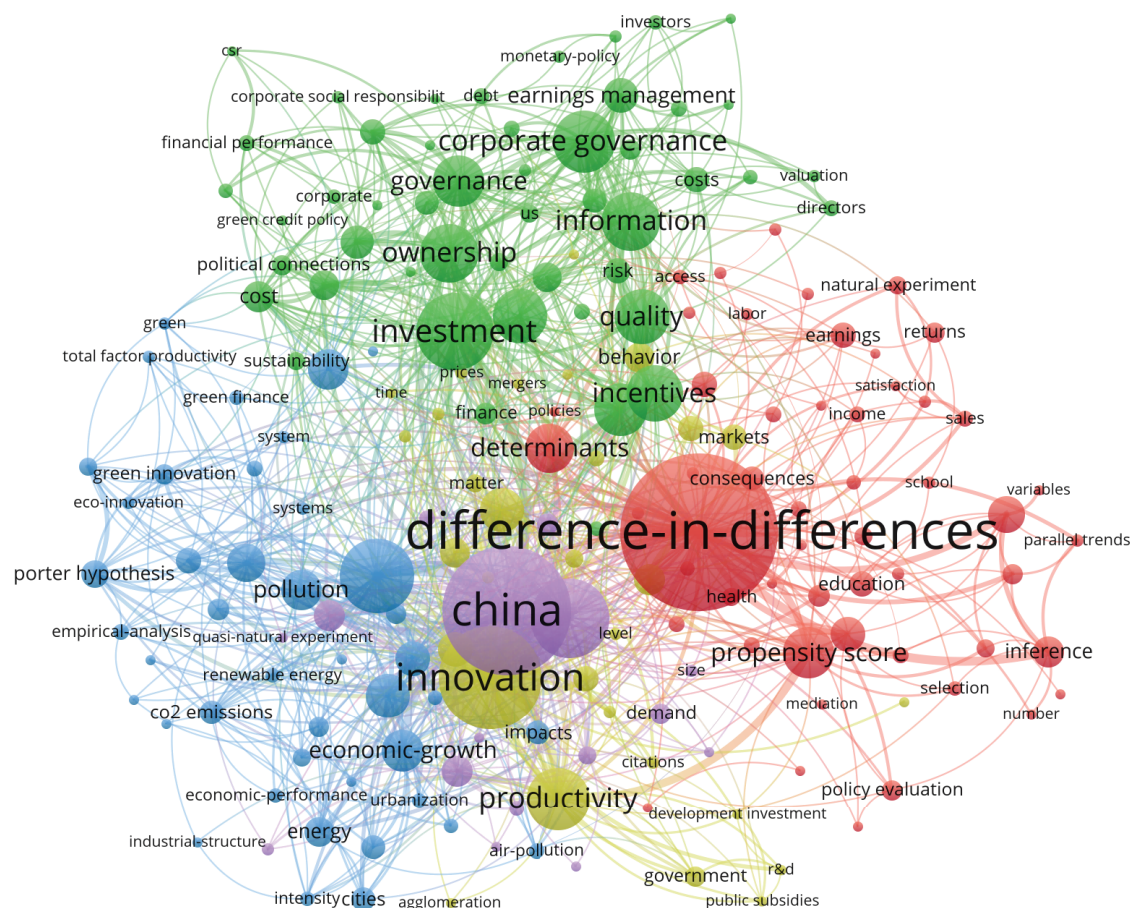
rise of high-dimensional data environments—where researchers can access hundreds of covariates from administrative records, digital platforms, or satellite imagery—has exposed the limitations of traditional causal tools.[4] Conventional DID approaches, often implemented via linear fixed effects models, may falter in the presence of nonlinear confounding, covariate imbalance, or heterogeneous treatment response.[5] In such settings, machine learning (ML) offers promising tools to flexibly model complex data structures, select relevant variables, and improve predictive accuracy.[6] Yet, ML methods themselves are often prediction-oriented and lack the structural discipline needed for causal interpretation.[7] Bridging the divide between structure-driven identification and data-driven flexibility has become a central challenge for contemporary empirical research.[8] In response, a growing literature—most notably Double/Debiased Machine Learning (DDML)—has attempted to integrate the strengths of both paradigms.[9] This paper situates itself within this growing effort to reconcile causal structure and machine learning innovation, focusing particularly on the evolution of DID and its modern extensions.[10] Through a comprehensive review and methodological synthesis, we propose a new estimator—S-DIDML that structurally embeds DID within a residualized ML framework, enabling robust and interpretable causal inference in high-dimensional settings.[11]

## 2.Literature Mapping and Bibliometric Review

### 2.1 The Evolution and Diversification of Difference-in-Differences in Social Science Research
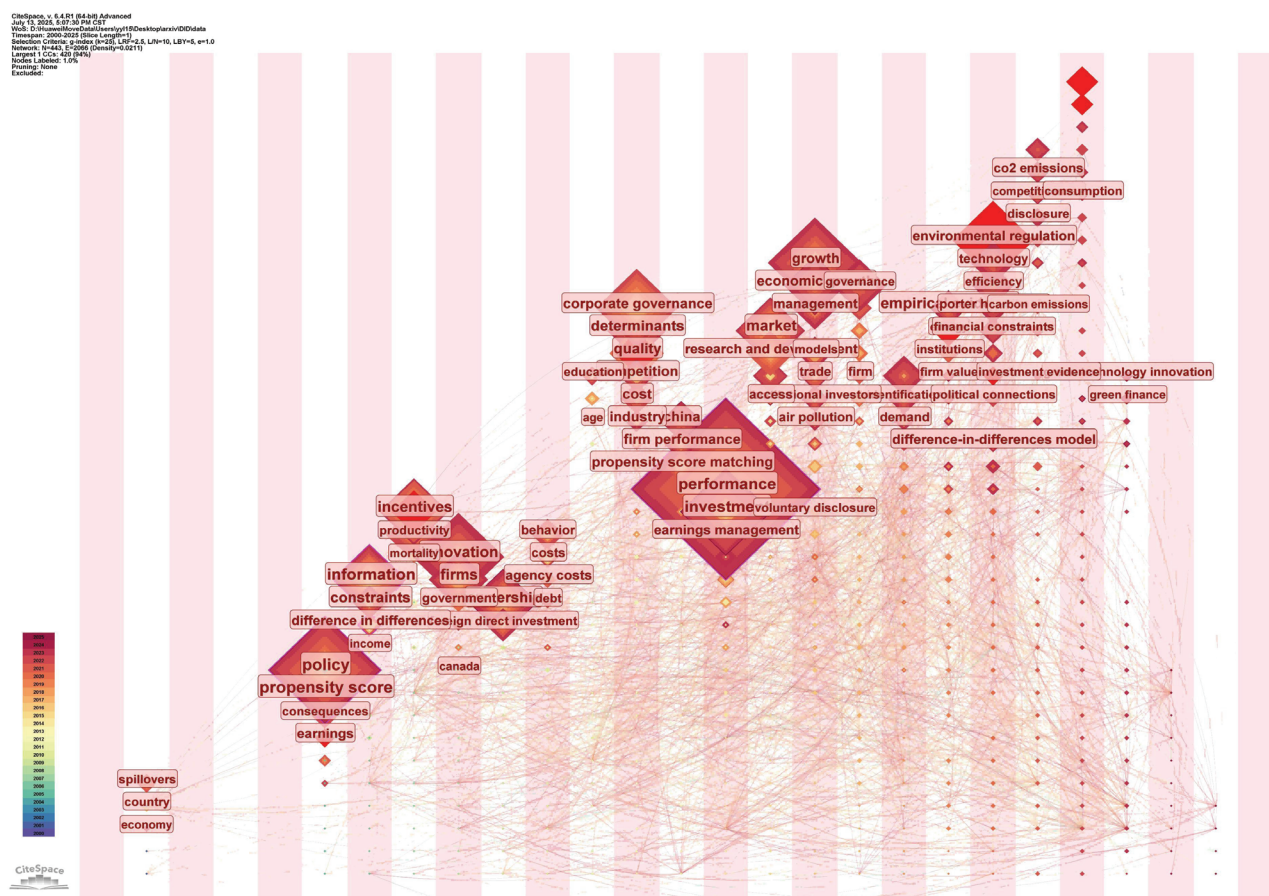
Over the past three decades, the Difference-in-Differences (DID) approach has evolved from a relatively simple econometric tool to a cornerstone methodology for causal inference in applied social sciences.[12] First widely adopted in labor economics and public policy evaluation (Ashenfelter & Card, 1985), DID methods are now routinely employed to assess the causal impacts of reforms and shocks across domains such as education, environmental regulation, corporate governance, and innovation policy.[13] The conceptual strength of DID lies in its structural identification strategy, which exploits longitudinal



*Figure 1 Co-occurrence clustering map of keywords in DID-related research*

variation between treatment and control units, assuming parallel trends in the absence of treatment.[14] Our bibliometric analysis of 500 articles indexed in Web of Science (2000–2024) reveals an exponential growth in DID-related publications, especially after 2010. This growth reflects the method's diffusion beyond economics into interdisciplinary domains, including political science, environmental studies, and development research. The VOSviewer keyword co-occurrence map (Figure 1) illustrates how DID research has clustered around several thematic poles: green and environmental economics, anchored by terms such as pollution, green finance, and Porter hypothesis; corporate governance and financial policy, focused on ownership, incentives, and earnings management; innovation and productivity analysis, with keywords like urbanization, China, and total factor productivity; and public policy and human capital evaluation, integrating DID with propensity score methods, education, and healthcare outcomes. The temporal cluster visualization from CiteSpace (Figure 2) suggests a clear trajectory of conceptual evolution.[15] Early DID studies (pre-2005) were tightly linked to labor market outcomes and macroeconomic shocks.[16] However, since 2015, the frontier has shifted toward applications in environmental regulation, technological innovation, and firm behavior under asymmetric information, often in the context of emerging markets such as China.[17] These newer applications frequently involve complex multi-treatment environments, time-varying exposures, and heterogeneous policy responses across sectors—conditions that test the limits of classical DID assumptions. Methodologically, DID estimation has undergone substantial refinement.[18] Influential work by Abadie (2005) extended DID into the semiparametric domain, allowing for richer covariate structures via matching.[19] More recently, Callaway and Sant'Anna (2021) and Sun and Abraham (2021) introduced robust estimators for staggered adoption designs, correcting for treatment timing heterogeneity that plagues two-way fixed effects models.[20] These innovations address several longstanding concerns: violation of parallel trends across heterogeneous units; bias under dynamic treatment effects and anticipation; and negative weighting in TWFE regression models. Despite these advances, traditional DID remains fundamentally limited in high-dimensional settings. Most DID implementations rely on linearity assumptions and low-dimensional covariate control, which become fragile when dealing with hundreds of potential confounders from administrative records, digital platforms,

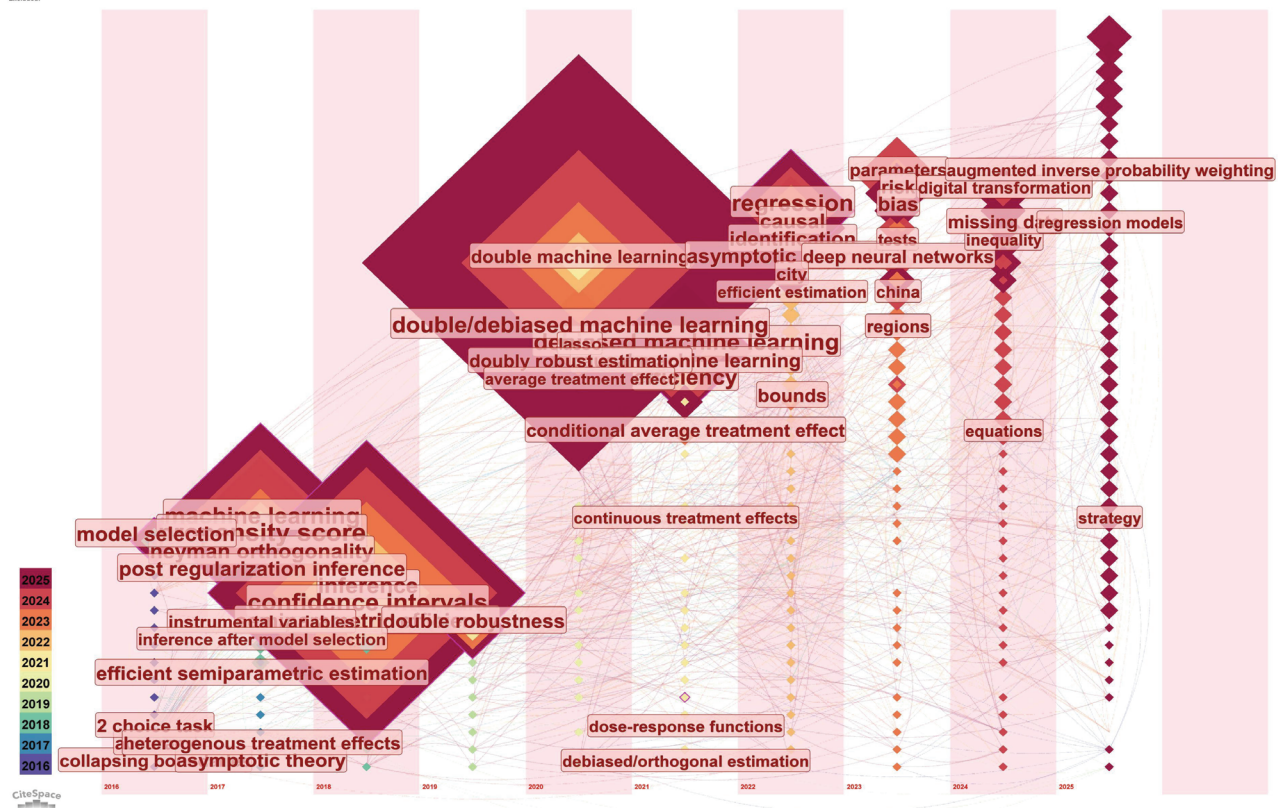*Figure 2 Temporal evolution diagram of keywords in DID-related research*

or geospatial data. Additionally, classical DID is typically estimated at the average treatment effect (ATE) level, with little consideration for treatment effect heterogeneity (HTE) across firms, regions, or demographic groups. Recent attempts to augment DID with machine learning techniques—particularly in pre-processing stages using propensity score estimation, kernel matching, or covariate balancing—remain largely modular rather than integrative.[21] They seldom modify the core identification strategy or link to formal orthogonalization procedures, as found in Double Machine Learning (DML) frameworks.[22] In summary, DID has achieved a remarkable degree of empirical relevance and institutional acceptance across disciplines.[23] However, the emerging landscape of complex, high-dimensional policy environments exposes critical limitations in classical DID frameworks.[24] To remain a credible tool for modern causal analysis, DID methods must evolve in two key directions: (1) integration with flexible, data-adaptive modeling approaches capable of handling nonlinearities and latent confounding, and (2) preservation of structural interpretability, ensuring that causal parameters remain transparent and policy-actionable.[25] This methodological tension motivates the development of S-DIDML, a new estimator introduced in this study that embeds DID within a double residualized machine learning framework.[26] S-DIDML retains the core logic of temporal comparison while incorporating high-dimensional learning and orthogonalized estimation.[27] As we will show, this approach offers a principled and scalable path toward heterogeneity-aware, high-dimensional causal inference grounded in quasi-experimental logic.[28]

## 2.2 The Rise of Double/Debiased Machine Learning: From Orthogonalization to High-Dimensional Causal Inference

In response to the growing need for reliable causal inference in high-dimensional settings, the past decade has witnessed the rapid rise of Double Machine Learning (DML) and Debiased Machine Learning (DDML) frameworks.[29] Rooted in the econometric concept of orthogonalized moment equations and powered by modern machine learning-based nuisance parameter estimation, these methods have substantially expanded the empirical frontier for researchers dealing with complex treatment assignment, heterogeneity, and large-scale observational data.[30] DML and its variants provide a principled way to separate the prediction task (nuisance estimation) from the estimation of causal effects, by leveraging Neyman orthogonality, sample-splitting, and cross-fitting to obtain asymptotically valid estimators even in the presence of flexible, nonparametric models.[31] The bibliometric analysis of 178 core papers indexed in Web of Science (2016–2024) confirms the explosive growth of this literature.[32] The VOSviewer co-occurrence network (Figure 3) centers around key themes including causal inference, Neyman orthogonality, cross-fitting, efficiency, and heterogeneous treatment effects, all closely linked with practical implementation strategies such as lasso, random forest, double robustness, and partially linear models.[33] Meanwhile, the CiteSpace temporal map (Figure 4) highlights a condensed knowledge burst in 2019–2022, during which seminal works formalized inference after machine learning (Chernozhukov et al.[34], 2018), introduced double/debiased orthogonal estimation (Newey & Robins, 2018), and developed scalable estimators for average and conditional treatment effects (Athey & Wager, 2019; Farrell et al., 2021). These methods have now been codified into a general paradigm: use ML to estimate nuisance components, and then debias the effect estimator via orthogonalization. At the theoretical level, DDML builds on semiparametric efficiency theory, combining flexible first-stage ML tools (e.[35]g., penalized regression, tree-based ensembles, deep nets) with second-stage doubly robust score functions. Estimators are designed to satisfy conditions such as local robustness, asymptotic linearity, and root-n convergence, making them particularly suited for policy evaluation under approximate sparsity or nonlinear selection. A critical innovation is the use of cross-fitting to mitigate overfitting bias in high-capacity learners, enabling the use of complex models like neural networks or gradient boosting within a valid inferential framework. In empirical applications, DDML has rapidly diffused into fields such as health economics, education policy, taxation, and labor market discrimination, often in the context of heterogeneous treatment effects or partial identification.[36] For instance, Athey, Tibshirani, and Wager (2019) propose Causal Forests to estimate treatment effects conditional on covariates, while Chernozhukov et al.[37] (2020) extend DDML to quantile regression and instrumental variables. These contributions have redefined the scope of credible causal inference, enabling researchers to shift from estimating a single average treatment effect toward recovering rich heterogeneity structures and distributional effects, all while preserving valid statistical inference. Nevertheless, current DDML applications still face several constraints.[38] First, many frameworks are

*Figure 3 Keyword co-occurrence clustering map of DDML-related research*



*Figure 4 Temporal evolution diagram of keywords related to DDML research*

cross-sectional or static, lacking explicit integration with panel designs, staggered treatments, or temporal structures typical of Difference-in-Differences research.[39] Second, the adoption of ML within causal inference remains primarily focused on prediction-quality improvements rather than structural causal modeling.[40] There is limited attention to embedding domain-specific identification strategies (e.g., timing, policy eligibility rules, group heterogeneity) within the estimation process. Finally, while DML allows for flexible controls, it does not by itself address the interpretability challenge: the black-box nature of some ML learners can obscure the causal estimand and undermine transparency in applied policy contexts. These challenges underscore the necessity of an integrative approach—one that unifies the structure-driven identification logic of quasi-experimental designs with the data-adaptive capacity of machine learning. The S-DIDML framework proposed in this paper is precisely such an effort. By embedding the Difference-in-Differences design into a DDML architecture, we enable researchers to retain structural interpretability while achieving statistical robustness in high-dimensional, dynamic, and heterogeneous settings.

## 2.3 Structural Causal Inference Meets Machine Learning: Toward a Unified Framework for Transparent and Scalable Policy Evaluation

The increasing availability of high-dimensional observational data has amplified the need for causal inference methods that are both structurally interpretable and computationally scalable. This has led to a surge in research endeavors at the intersection of structural causal inference and machine learning (ML), with a focus on preserving the clarity of model-based identification strategies while leveraging the flexibility and generalization capabilities of ML algorithms. Conventional econometric models rely heavily on parametric assumptions, while pure machine learning (ML) approaches optimize prediction without providing inference guarantees. In contrast, this hybrid literature aims to formalize data-adaptive structural estimators that are consistent, efficient, and interpretable for policy evaluation. A comprehensive review of 62 high-impact publications (2015–2024) from Web of Science reveals that this integration effort is centered on several converging themes. As demonstrated in the VOSviewer map (Figure 5), the predominant themes encompass causal inference, propensity score matching, double robust estimation, instrumental variables, and generalized random forests, frequently intersecting with applied domains such as healthcare, education, energy policy, and digital platforms. In terms of methodology, big data analytics, deep learning, bias reduction, and heterogeneous treatment effects are frequently mentioned alongside fundamental statistical terms such as efficiency, identification, and confidence intervals. Concurrently, the CiteSpace timeline (Figure 6) underscores three significant phases in the field's intellectual trajectory: an initial predilection for efficient semiparametric estimation (2015–2018), a transition to policy-aware machine learning models (2019–2021), and a recent proliferation of causal machine learning for complex treatment regimes (2022–2024). A prominent line of research focuses on integrating graphical causal models or potential outcomes frameworks into machine learning (ML) pipelines. For instance, Louizos et al. (2017) introduced Causal Effect Variational Autoencoders, which utilize latent variable modeling to disentangle treatment effects from confounding. In their 2017 publication, Shalit et al. put forward a proposal for the use of Counterfactual Regression Networks in the context of learning individualized treatment effects under the condition of a covariate shift. These architectures signify a substantial departure from the conventional black-box prediction paradigm, marking a transition toward model-aware architectures that demonstrate a heightened level of respect for treatment assignment mechanisms. Concurrently, within the econometrics tradition, researchers such as Pearl (2009) and Bareinboim & Pearl (2016) advocate for structural causal diagrams (SCMs) that formalize identification via do-calculus. These diagrams provide a framework for integrating data with domain knowledge to estimate causal parameters under selection bias or mediation. The synthesis of these approaches has resulted in the development of highly effective new estimators. It is noteworthy that generalized random forests (Athey et al., 2019) facilitate the estimation of conditional average treatment effects (CATE) through a partition-based approach. Additionally, causal boosting and meta-learners (e.g., T-learner, X-learner) adapt conventional ensemble learners to optimize policy evaluation tasks. These innovations frequently demonstrate superior performance in terms of bias-variance tradeoffs when confronted with classical regression models, particularly in scenarios where treatment effects exhibit significant heterogeneity and intricate interactions are prevalent. However, this emerging body of literature also reveals significant tensions between structural rigor and algorithmic complexity. First, while machine learning (ML) tools excel at

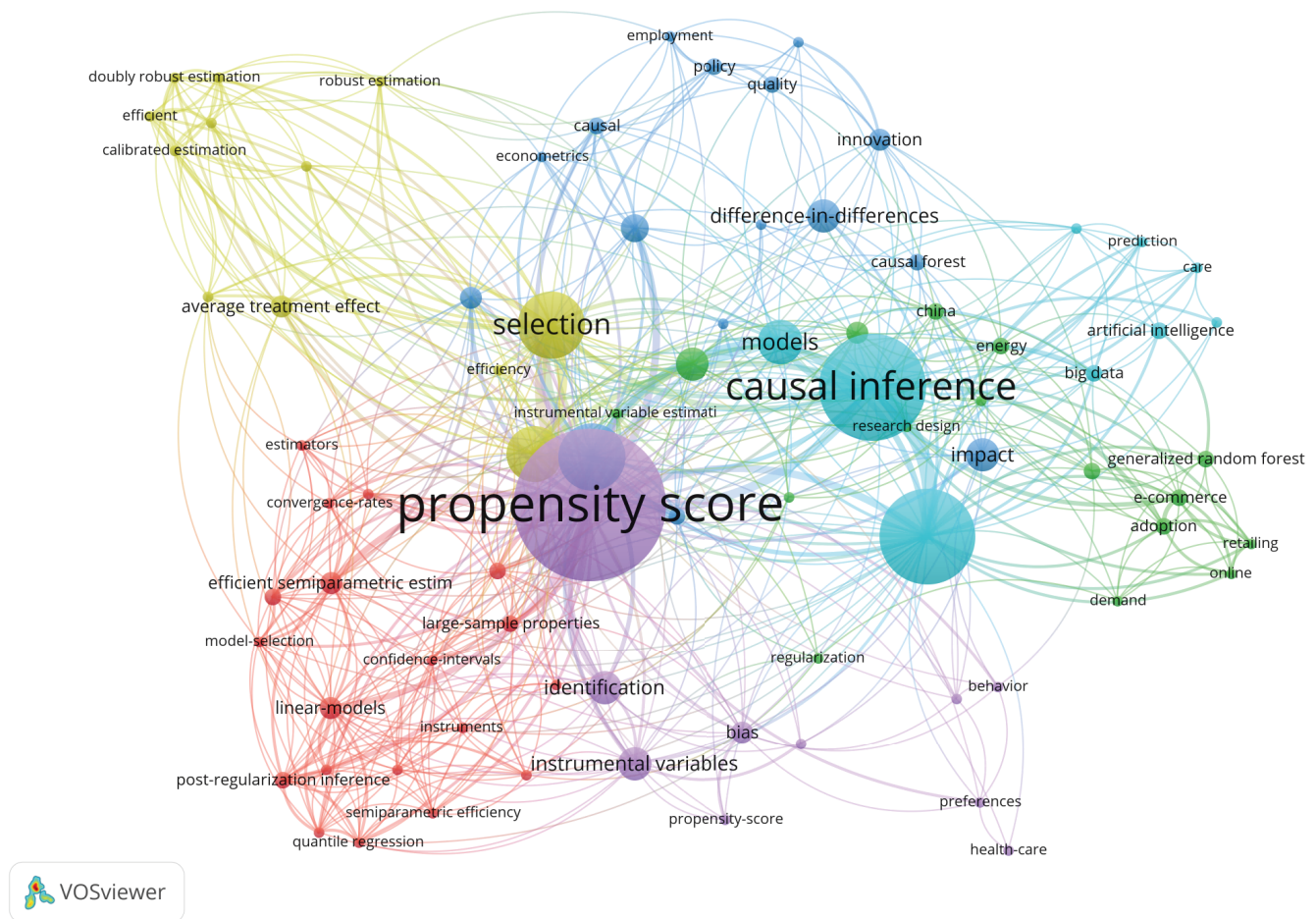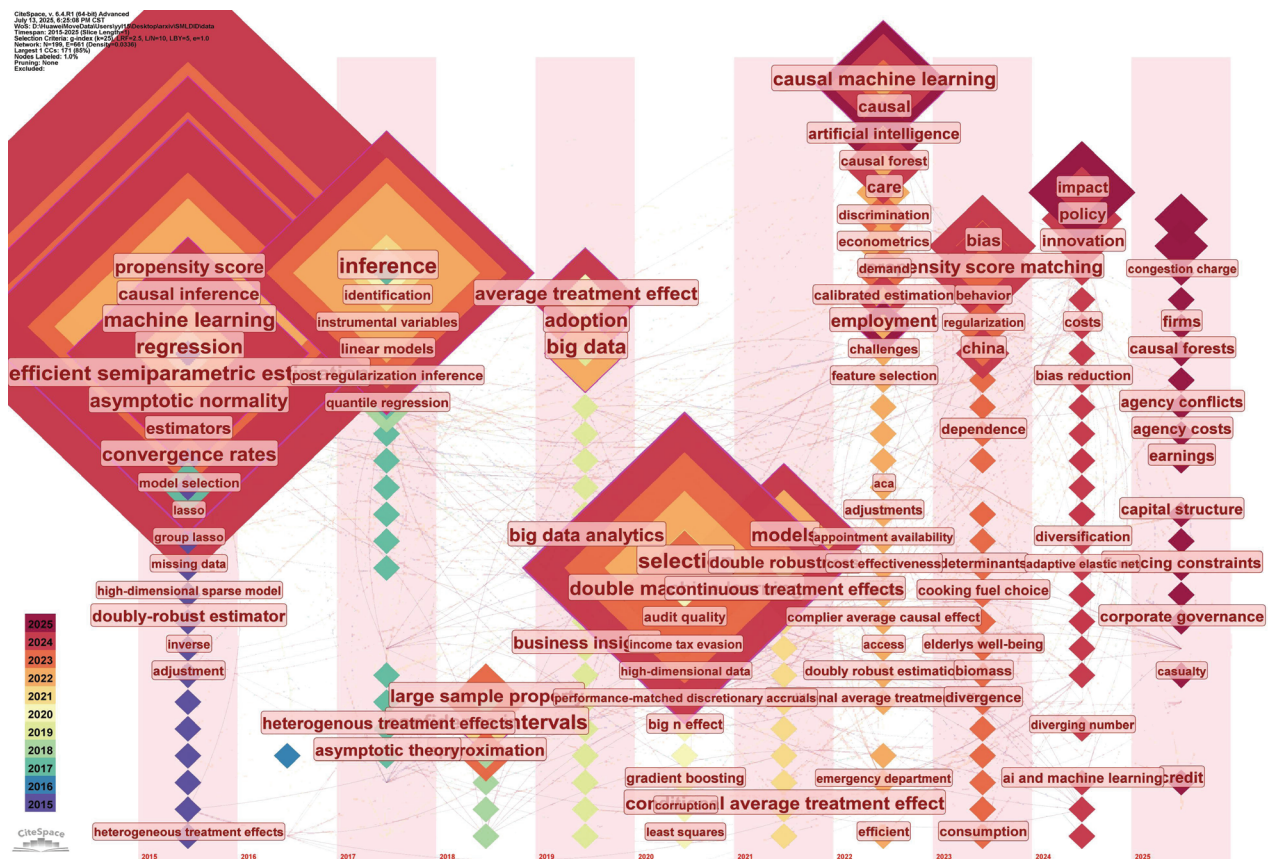*Figure 5DDML & DID Co-occurrence Cluster Map of Key Research Terms*



*Figure 6 Keyword Temporal Evolution Diagram of DDML&DID Integration Research*

approximating nuisance functions, they often lack explicit representation of causal assumptions. This hinders the evaluation of identification strength and the falsification of key assumptions (e.g., unconfoundedness). Secondly, structural models offer transparency and counterfactual interpretability; however, they are computationally fragile in high-dimensional settings, particularly under limited overlap or weak instruments. The absence of a unifying inferential theory further complicates the integration of frequentist inference with black-box learners, particularly in the context of interpreting uncertainty surrounding causal estimates. The proposed S-DIDML framework in this paper addresses these gaps by integrating difference-in-differences identification logic, double machine learning estimation, and domain-aware structural assumptions into a unified design. The model maintains clear identification through the use of parallel trends logic and temporal contrasts, enhancing estimation robustness using orthogonalized residualization and cross-fitting. Additionally, it enables high-dimensional adjustment through the implementation of supervised machine learning. Consequently, S-DIDML makes a significant contribution to the extant literature by offering a transparent, scalable, and statistically valid framework for policy evaluation with heterogeneous effects and complex treatments.

# 3.Gaps and Research Needs

A brief commentary is provided in the following table 1.

*Table 1 Comparison Table of Characteristics of Different Methods*

| Method category | Core limitations | Performance | References |
|---|---|---|---|
| Machine learning(ML) | Lack of interpretability and structural modeling capability | Black-box models are difficult to interpret causal pathways; lack economic behavioral constraints or theoretical assumptions; and are non-analyzable for decision-making processes. | Breiman (2001); Chernozhukov et al. (2018); Molnar (2022) |
| Difference-in-Differences (DID) | Unable to handle high-dimensional covariates and complex heterogeneity structures | High-dimensional conditions lead to model instability; Assumptions (such as parallel trends) are difficult to verify; Unable to naturally identify multi-group, multi-time-point heterogeneous effects | Goodman-Bacon (2021); Roth et al. (2022) |
| Double/Debiased ML (DML/ DDML) | Lack of structural embedding and multi-period adaptability | Assuming cross-sectional data is being processed, it is difficult to handle multi-period policies and temporal dynamics; divorced from the context of economic theory, the explanatory power is weak. | Chernozhukov et al. (2018); Imai & Kim (2021); Kennedy (2022) |
| HTE Estimation | Structural interpretation is lacking, and the mechanism transparency is low. | Causal Forests / Meta-Learners can detect heterogeneity but cannot explain the source of differences; they are prone to producing a "mechanism void" in policy interpretation bias. | Wager & Athey (2018); Künzel et al. (2019); Heckman & Vytlacil (2005) |

## 3.1Machine Learning in Causal Inference: Challenges of Interpretability and Weak Structural Assumptions

Machine learning (ML) techniques have exhibited remarkable success in high-dimensional prediction tasks. However, their application in causal inference is encumbered by fundamental limitations, particularly with regard to interpretability and structural grounding. This tension stems from the fact that the majority of supervised machine learning (ML) algorithms are designed to maximize predictive accuracy rather than to ensure causal identifiability. This phenomenon, as famously articulated by Pearl (2009), is often referred to as "the algorithmization of association, not causation." A significant challenge is that machine learning (ML) models generally operate under agnostic data-generating assumptions, exhibiting a lack of

a clear mapping to domain-relevant structural models. As Athey and Imbens (2019) emphasize, "machine learning excels at reducing prediction error but often leaves causal structure unspecified," which means treatment effects may be detected without being explainable in terms of mechanisms, counterfactual logic, or institutional design. To illustrate, while tree-based learners, such as random forests or gradient boosting, are capable of detecting treatment heterogeneity, they lack the capacity to discern whether this heterogeneity is attributable to observed policy variation, unmeasured confounding, or spurious interactions. The absence of structural constraints in machine learning (ML) introduces a secondary risk, which is the impediment of credible identification of causal parameters. In contrast to parametric econometric models, which explicitly encode assumptions such as exclusion restrictions, monotonicity, or sequential ignorability, machine learning (ML) models frequently estimate flexible functions without the use of guiding restrictions. As Louizos et al. (2017) have noted, this can result in high-variance estimators or misleading results, particularly when treatment assignment is non-random or when confounding is only partially observed. Furthermore, when machine learning (ML) is applied in a naive manner in causal contexts, it frequently violates the orthogonality conditions necessary for valid inference, unless careful debiasing or orthogonal score construction is employed (Chernozhukov et al., 2018). Interpretability is a closely related concern. As Kitson (2025) emphasizes, while techniques such as SHAP or LIME offer local approximations for black-box models, they "do not constitute structural explanations of the data-generating process" and cannot substitute for a formal causal model. The dearth of counterfactual semantics and policy-relevant parameters in standard machine learning (ML) frameworks engenders challenges in substantiating findings for regulatory, legal, and institutional decision-making. Collectively, these limitations point to a broader epistemological issue: without structural assumptions, causal statements derived from machine learning (ML) risk being descriptive rather than explanatory. This weakens the scientific value of such analyses, particularly in social science domains that rely heavily on theoretical grounding, historical context, and institutional realism. In order to address these concerns, the extant literature proposes a combination of machine learning (ML) with formal structural inference frameworks, such as Structural Equation Models (SEMs) (Pearl, 2009), structural score functions (Newey & Robins, 2018), or graphical causal modeling (Bareinboim & Pearl, 2016). However, these integrations remain in their infancy and have yet to be widely adopted in empirical work, particularly in policy evaluation with staggered or longitudinal designs.

## 3.2 Limitations of Classical Difference-in-Differences: Covariate Dimensionality and Identification Fragility

The Difference-in-Differences (DID) framework has become a foundational method in the field of applied econometrics for estimating causal effects from observational panel data. The study's fundamental appeal stems from its use of an intuitive identification strategy, which involves the comparison of outcome trends between treated and control groups. This approach is predicated on the assumption of parallel trends in the absence of treatment. However, despite its widespread use, recent theoretical and empirical developments have highlighted critical limitations of DID, particularly in the face of high-dimensional covariates, treatment timing heterogeneity, and violations of baseline assumptions. A primary concern pertains to the management of covariates and the issue of model misspecification. In the field of data science and statistics, traditional DID (Diffusion Interrupted Difference) models are typically estimated using two-way fixed effects (TWFE) regression. In this estimation method, control variables are either excluded or incorporated linearly. In practice, this limitation restricts the capacity of DID to adjust for nonlinear, time-varying, or high-dimensional confounders—precisely the kinds of complexities that are prevalent in modern administrative, firm-level, or geospatial data. As Abadie (2005) observed, even moderate deviations from parametric assumptions in the outcome model can result in biased estimates. The magnitude of these risks is amplified under two conditions: first, when the covariate space expands, and second, when interactions between covariates and treatment status are neglected. Furthermore, as demonstrated by Goodman-Bacon (2021), the implementation of TWFE in staggered adoption settings can result in the negative weighting of treatment effects, thereby introducing bias and compromising the interpretability of causal relationships. This phenomenon occurs because the estimator aggregates over comparisons across different timing groups, some of which may act as implicit controls for others. In such cases, the estimate no longer represents a clear causal contrast between treated and untreated units, especially if treatment effects are dynamic or heterogeneous across cohorts. These insights have contributed to the development of more robust DID estimators, such

as those proposed by Callaway and Sant'Anna (2021), which explicitly account for variation in treatment timing and allow for group-time-specific treatment effects. A further challenge lies in the robustness to trend violations. Although DID is often justified by informal graphical checks or pre-trend tests, these approaches suffer from low power and subjective interpretation. Roth (2023) offers a critique of the overreliance on pre-trends as a robustness diagnostic and proposes formal methods to account for uncertainty in parallel trend assumptions. The author demonstrates that even mild violations can substantially alter inference. Furthermore, in numerous policy contexts—including rolling interventions or gradually implemented regulations—the concept of a singular pre-treatment trend is often deemed implausible, necessitating more adaptable and data-driven trend modeling methodologies. The conventional DID framework is inadequate in addressing high-dimensional settings, where the number of potential covariates exceeds the sample size. In such circumstances, linear fixed-effects models may become unstable or inapplicable, and post-hoc covariate balancing or matching procedures may fail due to poor overlap or extreme weights. Attempts to augment the difference-in-difference (DID) method with machine learning for pre-processing (e.g., via propensity score estimation) are often modular rather than integrative. This failure to embed machine learning (ML) into the estimation stage or maintain orthogonality necessary for valid inference has been observed in many cases. While DID remains a powerful identification tool, its classical implementations face substantial limitations in modern empirical settings characterized by high-dimensional data, staggered treatments, heterogeneous effects, and limited trend credibility. These limitations have motivated the development of structurally grounded, machine-learning-enhanced DID frameworks—such as the S-DIDML estimator proposed in this study—that preserve the logic of temporal comparison while enhancing robustness, flexibility, and theoretical coherence.

## 3.3 Double Machine Learning: Weak Integration with Economic Structure and Limited Temporal Flexibility

Double Machine Learning (DML) has emerged as a powerful framework for causal inference in high-dimensional settings. By leveraging orthogonal moment conditions, sample-splitting, and cross-fitting, DML enables consistent estimation of treatment effects even when nuisance functions (e.g., propensity scores or outcome regressions) are estimated using complex, flexible machine learning methods (Chernozhukov et al., 2018). However, despite its strong statistical foundation, DML remains structurally minimalistic and temporally constrained, making it ill-suited for many empirical contexts encountered in economic policy research. A salient limitation pertains to DML's predilection for cross-sectional or static treatment settings. The majority of DML implementations are developed under the assumption of a single treatment decision per unit. This limitation precludes the application of these methods to scenarios such as multi-period treatments, staggered adoption, or event-time heterogeneity, which are common in education, tax, environmental, or labor policy studies. As Imai and Kim (2021) have demonstrated, the application of DML in panel data contexts can result in biased estimates if the temporal structure of treatment assignment and potential outcomes is not adequately considered. This limitation restricts the applicability of standard DML in real-world settings where treatments unfold over time and responses are dynamic. Additionally, while DML offers flexibility in estimating nuisance functions, it does not inherently ensure structural interpretability. The majority of DML applications are implemented without integrating formal economic assumptions—such as rational behavior, instrumental monotonicity, policy discontinuities, or selection mechanisms—into the estimation architecture. Consequently, DML estimators have the capacity to discern local causal effects, yet they frequently lack a comprehensive explanation of their theoretical underpinnings or their economic ramifications. Heckman and Pinto (2019) contend that such "structure-free" approaches carry the risk of producing effects that are "statistically significant but policy-irrelevant," particularly in domains that necessitate behavioral modeling or institutional contextualization. The efficacy of DML's inference guarantees is contingent upon the fulfillment of specific orthogonality and rate conditions. However, these conditions may become ineffective in settings characterized by low overlap, intricate treatment interactions, or feedback loops over time. For instance, Kennedy (2022) demonstrates that conventional cross-fitting and orthogonalization procedures may exhibit suboptimal performance when the treatment assignment mechanism itself is endogenous to outcomes (e.g., performance-based subsidies or policy-induced behavioral changes). This further underscores the necessity for structural restrictions to ensure causal stability and policy extrapolation. Another practical challenge lies in embedding DML within quasi-experimental designs,

such as regression discontinuity, instrumental variables, or difference-in-differences frameworks. Recent endeavors have sought to expand the DML framework to encompass instrumental variables (Chernozhukov et al., 2022) and event studies (Roth and Sant'Anna, 2023). However, these approaches remain in their nascent stages and frequently exhibit the ambiguity and opacity characteristic of traditional structural designs. The tension between statistical optimality and domain-specific identification logic remains unresolved. Collectively, these limitations suggest that while DML is a powerful statistical tool, its effectiveness in economic and social science applications depends on its integration with structural frameworks that reflect timing, group heterogeneity, policy rules, and institutional mechanisms. In the absence of such integration, DML faces the potential to evolve into a technically sophisticated yet substantively opaque estimation strategy.

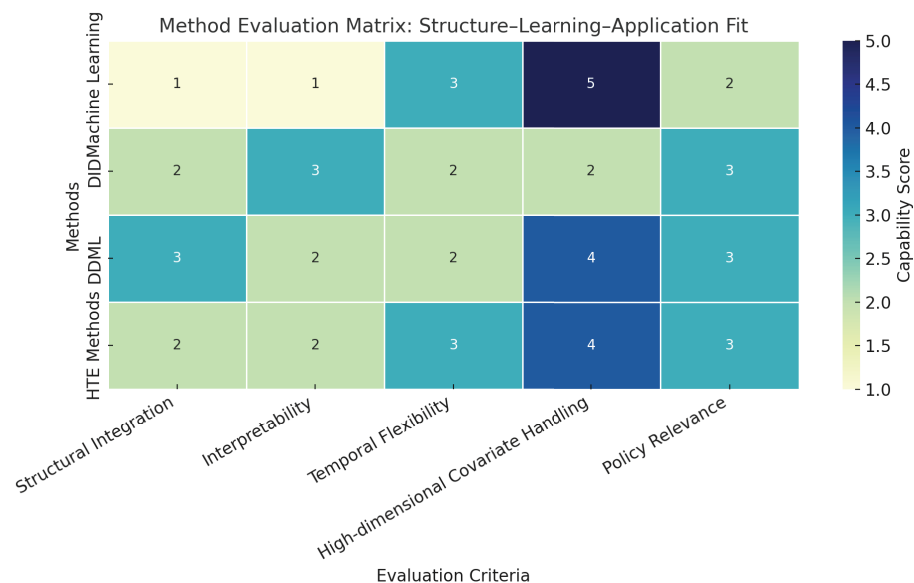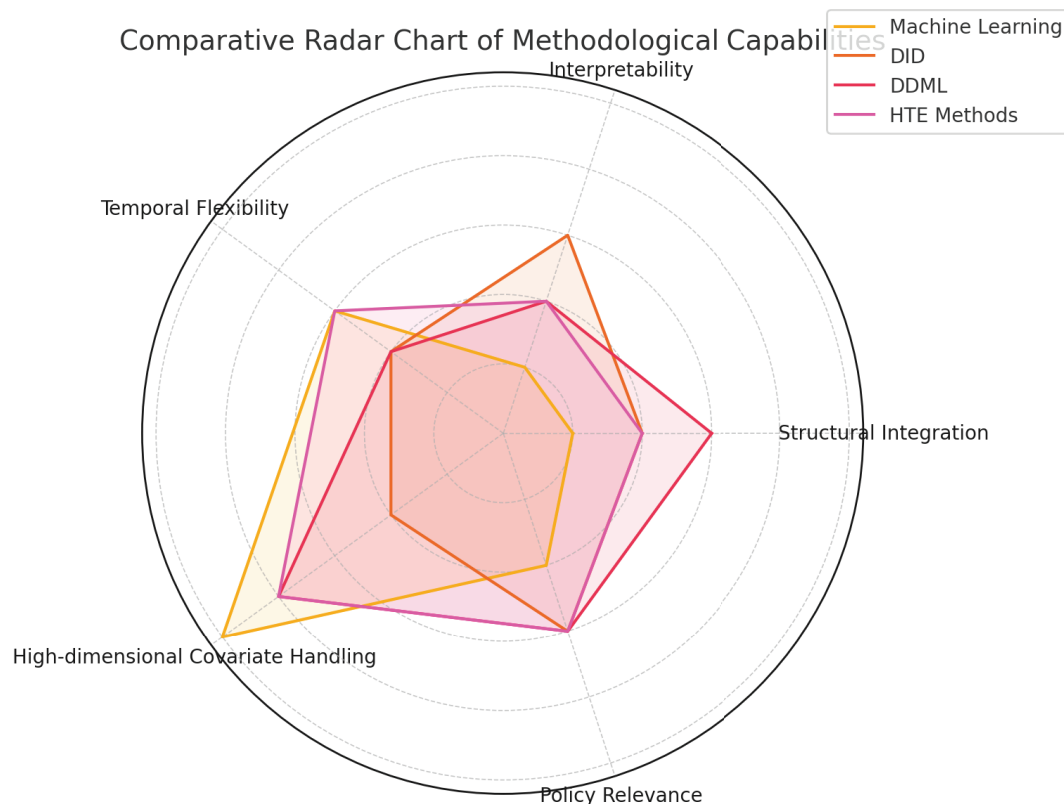*Figure 7 Matrix heatmap estimation diagram of characteristics for various research methods*



*Figure 8 Radar estimation diagram of the characteristics of various research methods*

# 4.The S-DIDML Framework

The S-DIDML framework is constructed as a five-step semiparametric estimation pipeline, designed to combine the temporal identification logic of Difference-in-Differences (DID) with the residualization and orthogonalization principles of Double Machine Learning (DML). This design enables robust causal inference in high-dimensional, staggered-treatment settings while preserving interpretability grounded in structural counterfactual logic.

Step 1: Policy Exposure Encoding and Panel Structuring

The initial step involves transforming raw data into a well-defined panel format. Treatment exposure variables $D_i t$ are encoded to reflect group-level policy adoption, including information on timing and staggered rollout. Cohort indicators $g(i)$ and time indicators t are constructed, enabling the identification of treatment dynamics across subgroups and time periods. This step reproduces the structural basis of traditional DID under staggered adoption scenarios.

Step 2: High-dimensional Nuisance Estimation via Machine Learning

To control for confounding in high-dimensional settings, flexible machine learning models are used to predict the outcome and treatment assignment based on covariates. Specifically, two nuisance functions are estimated:

$$g(X_i t) \approx E[Y_i t | X_i t], m(X_i t) \approx E[D_i t | X_i t]$$

These models are fitted using cross-fitting to ensure orthogonality and to mitigate overfitting. A wide range of supervised ML methods (e.g., random forests, boosting, neural networks) can be applied at this stage, depending on the structure of $X_i t$.

Step 3: Double Residualization of Outcome and Treatment Variables

Following estimation, the observed variables are residualized as follows:

$$\tilde{Y}_i t = Y_i t - \hat{g}(X_i t), \tilde{D}_i t = D_i t - \hat{m}(X_i t)$$

This double residualization process yields outcome and treatment variables that are orthogonal to the high-dimensional covariates $X_i t$, thereby satisfying the Neyman orthogonality condition necessary for valid second-stage inference.

Step 4: Structural DID Estimation on Residualized Quantities

The core causal effect is estimated by regressing the residualized outcome $\tilde{Y}_i t$ on the residualized treatment indicator $\tilde{D}_i t$, while incorporating group and time fixed effects. This regression preserves the cohort-time structure of DID and allows for dynamic, group-specific treatment effects:

$$\tilde{Y}_i t = \tau_g, t \cdot \tilde{D}_i t + \alpha_g + \lambda_t + \varepsilon_i t$$

Depending on the empirical setting, estimators such as Callaway–Sant'Anna (2021) or Sun & Abraham (2021) can be employed to estimate average or event-time treatment effects.

Step 5: Aggregation, Uncertainty Quantification, and Robustness Checks

Estimated group-time treatment effects $\tau_g$, t are aggregated into overall ATT or dynamic treatment effect curves. Standard errors are obtained using cross-fitting–compatible variance formulas or nonparametric bootstrap methods. Finally, robustness is assessed via falsification tests (e.g., placebo interventions), checking for pre-trend violations, and assessing overlap conditions.

# 5.Demonstrative Literature Applications: Where S-DIDML Can Be Used

The S-DIDML framework is not only theoretically robust but also practically versatile. Its design enables immediate adoption in several streams of empirical literature, especially where conventional DID or DDML frameworks face limitations due to high-dimensional covariates, staggered policy timing, or heterogeneity in treatment effects. Below, we outline four thematic domains where S-DIDML can provide substantial improvements in causal identification and inference quality.

## 5.1 Labor Economics: Evaluating Active Labor Market Policies (ALMPs)

Many ALMP evaluations rely on DID or event-study approaches (e.g., Kluve, 2010; Card et al., 2018), often using limited covariates due to multicollinearity concerns. However, modern administrative labor datasets now include thousands of features (firm size, tenure, dynamic local shocks). S-DIDML allows for robust estimation of heterogeneous effects of job subsidies or training programs across firms, sectors, or worker types, while maintaining structural interpretability of ATT.

## 5.2 Education Policy: School Reform, Curriculum Changes, and Tracking

Educational reforms (e.g., extending school years, STEM incentives, curriculum realignments) are often evaluated via DID with state or district fixed effects. However, treatment rollout is usually staggered, and student-level data are high-dimensional. By orthogonalizing outcomes with rich baseline test scores, socio-demographics, and parental inputs, S-DIDML enables dynamic treatment effect estimation at both cohort and demographic subgroup levels.

## 5.3 Fiscal and Tax Policy: Estimating Behavioral Responses with Administrative Tax Data

Tax policy changes (e.g., earned income tax credits, marginal rate changes) exhibit rich staggered designs but require flexible modeling of income dynamics, deductions, and family structure. Traditional regression-based DID models are poorly suited for such settings. S-DIDML accommodates high-dimensional pre-tax characteristics and allows for precise subgroup inference on labor supply elasticities or compliance behavior.
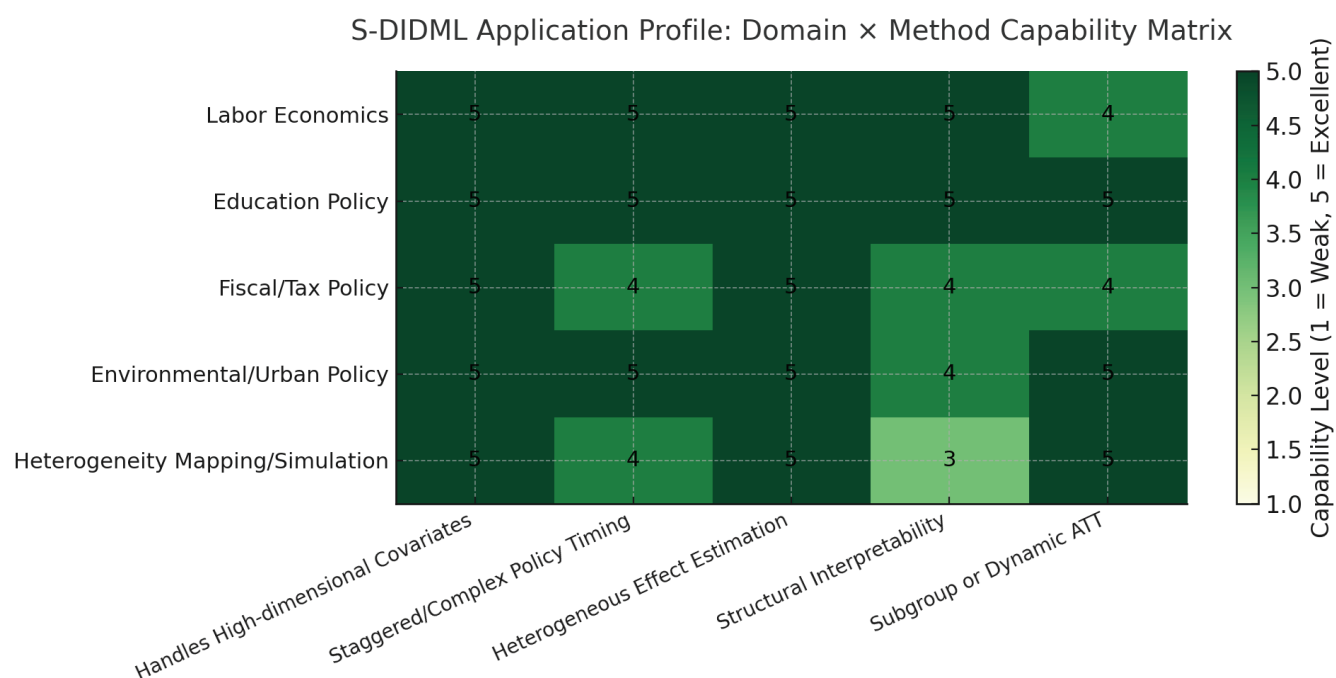
## 5.4 Environmental and Urban Policy: Evaluating Green Subsidies and Urban Interventions

Many environmental interventions (e.g., subsidies for electric vehicles, zoning regulations, pollution controls) are staggered and differ in intensity across space and time. Treatment heterogeneity is fundamental, and the policy environment is often rich in covariates (weather, geography, firm-level pollution histories). S-DIDML can flexibly estimate subgroup effects (e.g., by income decile or industry), while adjusting for spatial autocorrelation and nonlinearity.

## 5.5 Opportunities in Heterogeneity Mapping and Welfare Simulation

Beyond direct ATT estimation, S-DIDML can be embedded into policy simulation pipelines, enabling credible counterfactual mapping of treatment gains across the covariate space. For instance, it can assist in identifying which demographic groups benefit most from job guarantees or minimum wage hikes, using machine learning for heterogeneity partitioning while ensuring DID identification integrity.

*Figure 9 Matrix heatmap estimation diagram of S-DIDML framework applications across different domains*



## 6.Conclusion

This paper introduces the S-DIDML framework—a structural, semiparametric estimator designed to bridge the interpretability of Difference-in-Differences (DID) methods with the high-dimensional flexibility of Double Machine Learning (DML). In doing so, we aim to provide applied researchers in economics and the social sciences with a unified, theoretically grounded, and computationally feasible approach to estimating heterogeneous treatment effects in staggered policy contexts. We began by identifying a set of unresolved challenges in modern causal inference: the limited interpretability of pure machine learning estimators, the instability of conventional DID methods under high-dimensional controls, and the restricted scope of existing DDML estimators in handling multiple treatment periods or complex policy rollout. Through an extensive literature review

across DID, DDML, and structural ML, we demonstrated the methodological need for an integrated solution. To address these limitations, S-DIDML proposes a five-step estimation pipeline: (1) panel structuring and treatment timing encoding, (2) flexible nuisance estimation using ML, (3) double residualization for Neyman orthogonality, (4) structural DID regression for interpretable group-time effects, and (5) aggregation and robustness analysis. Each step is modular, theoretically motivated, and designed for transparency and scalability. We articulated the framework's principles—structural identification, orthogonality, and semiparametric adaptability—and illustrated its potential across key empirical domains including labor policy, education, taxation, and environmental regulation. Furthermore, we engaged critically with its current limitations, such as reliance on overlap, lack of interference modeling, and the need for unified subgroup inference. These issues mark important frontiers for future research in causal machine learning. S-DIDML is not intended to replace either traditional quasi-experimental methods or deep learning-based prediction tools. Instead, it acts as a conceptual and computational bridge: retaining causal interpretability grounded in economic theory, while leveraging the modeling capacity of modern ML for complex data environments. As empirical researchers face ever-larger datasets and increasingly heterogeneous policy designs, such hybrid frameworks are crucial for producing credible, robust, and policy-relevant insights. We envision S-DIDML not as a fixed model but as a flexible blueprint—one that invites further theoretical refinement, software development, and empirical adaptation. Its goal is not merely to improve estimation, but to foster a new generation of structurally informed, statistically rigorous causal inference in the high-dimensional era.

## Funding

## Conflict of Interests

The authors declare that there is no conflict of interest regarding the publication of this paper.

## Reference

[1] Abadie, A. (2021). Using synthetic controls: Feasibility, data requirements, and methodological aspects. Journal of Economic Literature, 59(2), 391–425. https://doi.org/10.1257/jel.20201405

[2] Athey, S., & Imbens, G. (2017). The state of applied econometrics: Causality and policy evaluation. Journal of Economic Perspectives, 31(2), 3–32. https://doi.org/10.1257/jep.31.2.3

[3] Callaway, B., & Sant'Anna, P. H. C. (2021). Difference-in-differences with multiple time periods. Journal of Econometrics, 225(2), 200–230. https://doi.org/10.1016/j.jeconom.2020.12.001

[4] Sant'Anna, P. H. C., & Zhao, J. (2020). Doubly robust difference-in-differences estimators. Journal of Econometrics, 219(1), 101–122. https://doi.org/10.1016/j.jeconom.2020.06.003

[5] Sun, L., & Abraham, S. (2021). Estimating dynamic treatment effects in event studies with heterogeneous treatment effects. Journal of Econometrics, 225(2), 175–199. https://doi.org/10.1016/j.jeconom.2020.09.006

[6] Chernozhukov, V., et al. (2018). Double/debiased machine learning for treatment and structural parameters. The Econometrics Journal, 21(1), C1–C68. https://doi.org/10.1111/ectj.12097

[7] Kennedy, E. H. (2022). Semiparametric theory and empirical processes in causal inference. Annual Review of Statistics and Its Application, 9, 151–176. https://doi.org/10.1146/annurev-statistics-040220-112545

[8] Nie, X., & Wager, S. (2021). Quasi-oracle estimation of heterogeneous treatment effects. Biometrika, 108(2), 299–319. https://doi.org/10.1093/biomet/asaa076

[9] Imai, K., & Kim, I. S. (2021). When should we use unit fixed effects regression models for causal inference with longitudinal data? American Journal of Political Science, 65(2), 448–466. https://doi.org/10.1111/ajps.12523

[10] Roth, J. (2023). Pre-test with care: How to test for parallel trends with multiple groups. Review of Economics and Statistics. https://doi.org/10.1162/rest_a_01207

[11] Angrist, J. D., & Pischke, J.-S. (2009). Mostly harmless econometrics: An empiricist's companion. Princeton University

Press.

[12] Wooldridge, J. M. (2021). Introductory Econometrics: A Modern Approach (7th ed.). Cengage Learning.

[13] Belloni, A., Chernozhukov, V., & Hansen, C. (2014). Inference on treatment effects after selection among high-dimensional controls. Review of Economic Studies, 81(2), 608–650. https://doi.org/10.1093/restud/rdt044

[14] Doudchenko, N., & Imbens, G. (2016). Balancing, regression, difference-in-differences and synthetic control methods: A synthesis. NBER Working Paper No. 22791.

[15] Borusyak, K., Jaravel, X., & Spiess, J. (2023). Revisiting event study designs. Econometrica, 91(1), 65–95. https://doi.org/10.3982/ecta20695

[16] Callaway, B., Goodman-Bacon, A., & Sant'Anna, P. H. C. (2023). Difference-in-differences with a continuous treatment. Journal of Econometrics. https://doi.org/10.1016/j.jeconom.2023.105417

[17] Athey, S., Tibshirani, J., & Wager, S. (2019). Generalized random forests. Annals of Statistics, 47(2), 1148–1178. https://doi.org/10.1214/18-AOS1709

[18] Künzel, S. R., Sekhon, J. S., Bickel, P. J., & Yu, B. (2019). Metalearners for estimating heterogeneous treatment effects using machine learning. PNAS, 116(10), 4156–4165. https://doi.org/10.1073/pnas.1804597116

[19] Hill, J. (2011). Bayesian nonparametric modeling for causal inference. Journal of Computational and Graphical Statistics, 20(1), 217–240. https://doi.org/10.1198/jcgs.2010.08162

[20] Imbens, G. W., & Rubin, D. B. (2015). Causal inference for statistics, social, and biomedical sciences: An introduction. Cambridge University Press.

[21] Roth, J., Sant'Anna, P. H. C., Bilinski, A., & Poe, J. (2022). What's trending in difference-in-differences? NBER Working Paper No. 31506.

[22] Duflo, E., Glennerster, R., & Kremer, M. (2008). Using randomization in development economics research. Handbook of Development Economics, 4, 3895–3962.

[23] Hastie, T., Tibshirani, R., & Friedman, J. (2009). The elements of statistical learning (2nd ed.). Springer.

[24] Oprescu, M., & Zhu, Y. (2023). Selective machine learning for heterogeneous treatment effect estimation. Journal of Causal Inference, 11(1). https://doi.org/10.1515/jci-2022-0021

[25] Wager, S., & Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. Journal of the American Statistical Association, 113(523), 1228–1242. https://doi.org/10.1080/01621459.2017.1319839

[26] Schmidheiny, K., & Siegloch, S. (2019). On event studies and distributed-lags in two-way fixed effects models. IZA Discussion Paper No. 12088.

[27] de Chaisemartin, C., & D'Haultfœuille, X. (2020). Two-way fixed effects estimators with heterogeneous treatment effects. American Economic Review, 110(9), 2964–2996. https://doi.org/10.1257/aer.20181169

[28] Xu, Y. (2017). Generalized synthetic control method: Causal inference with interactive fixed effects models. Political Analysis, 25(1), 57–76. https://doi.org/10.1017/pan.2016.2

[29] Kasy, M., & Sautmann, A. (2021). Adaptive treatment assignment in experiments for policy choice. Econometrica, 89(1), 113–132. https://doi.org/10.3982/ECTA17443

[30] Ben-Michael, E., Feller, A., & Rothstein, J. (2021). The augmented synthetic control method. Journal of the American Statistical Association, 116(536), 1789–1803. https://doi.org/10.1080/01621459.2021.1929245

[31] Hazlett, C. (2020). Regression discontinuity and heteroskedasticity. Political Science Research and Methods, 8(3), 551–566.

[32] Varian, H. R. (2014). Big data: New tricks for econometrics. Journal of Economic Perspectives, 28(2), 3–28.

[33] Breiman, L. (2001). Statistical modeling: The two cultures. Statistical Science, 16(3), 199–231.

[34] Knaus, M. C., Lechner, M., & Strittmatter, A. (2021). Machine learning estimation of heterogeneous labor market impacts of COVID-19 policies. Labour Economics, 72, 102054. https://doi.org/10.1016/j.labeco.2021.102054

[35] Bryan, G., Karlan, D., & Nelson, S. (2021). Commitment devices. Annual Review of Economics, 13, 561–583. https://doi.org/10.1146/annurev-economics-082420-112136

[36] Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. Journal of Educational Psychology, 66(5), 688–701.

[37] Heckman, J. J., & Vytlacil, E. (2007). Econometric evaluation of social programs, part I: Causal models, structural models and econometric policy evaluation. Handbook of Econometrics, 6, 4779–4874.

[38] Pearl, J. (2009). Causality: Models, reasoning, and inference (2nd ed.). Cambridge University Press.

[39] Heckman, J. J., Pinto, R., & Savelyev, P. A. (2013). Understanding the mechanisms through which an influential early childhood program boosted adult outcomes. American Economic Review, 103(6), 2052–2086.

[40] Imbens, G. W. (2020). Potential outcome and directed acyclic graphs: An overview. AEA Papers and Proceedings, 110, 358–361. https://doi.org/10.1257/pandp.20201008