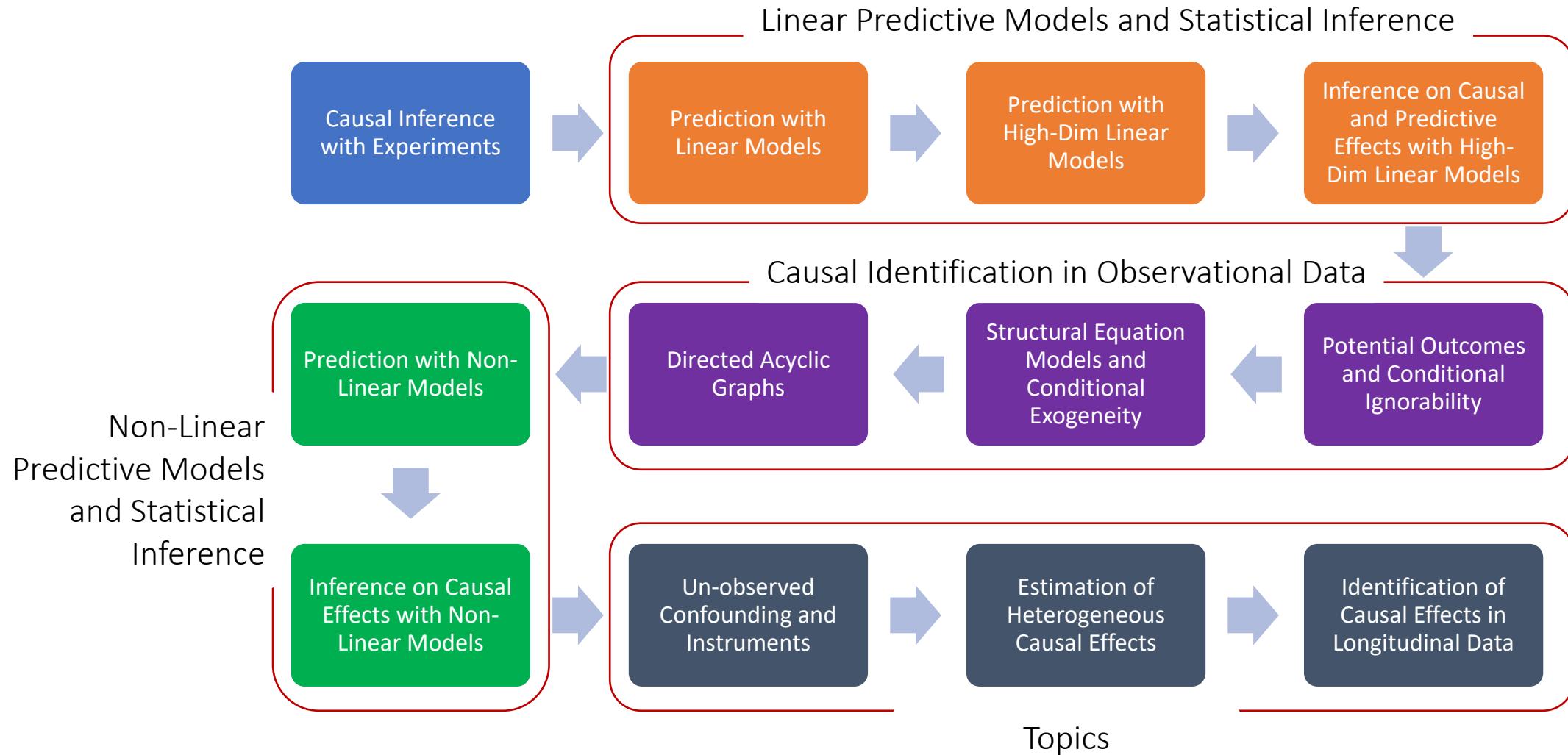
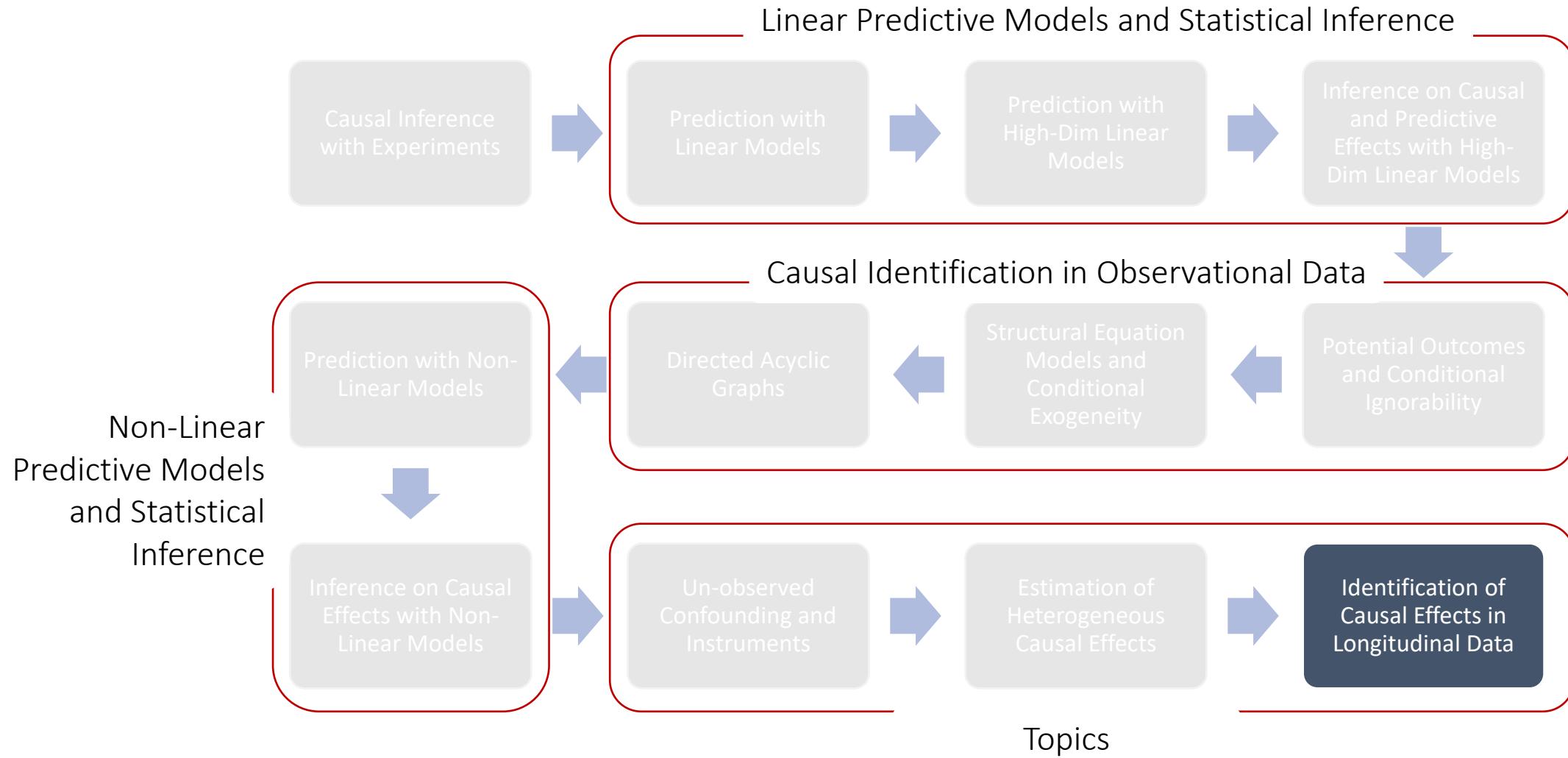


MS&E 228: Topics on Longitudinal Data and Causal Machine Learning

Vasilis Syrgkanis

MS&E, Stanford





Longitudinal Data

- Longitudinal data refer to datasets where we have multiple observations of the same unit over time

They present a separate set of new problems

- **Correlation:** samples stemming from each unit (or sometimes cluster of units from the same “site”) are correlated
- **Censoring/Missingness:** units typically drop out at (potentially not) random times before the end of the study
- **Dynamic treatments:** units are treated with multiple treatments over time in a manner that is adaptive and auto-correlated

Longitudinal Data

- Longitudinal data refer to datasets where we have multiple observations of the same unit over time

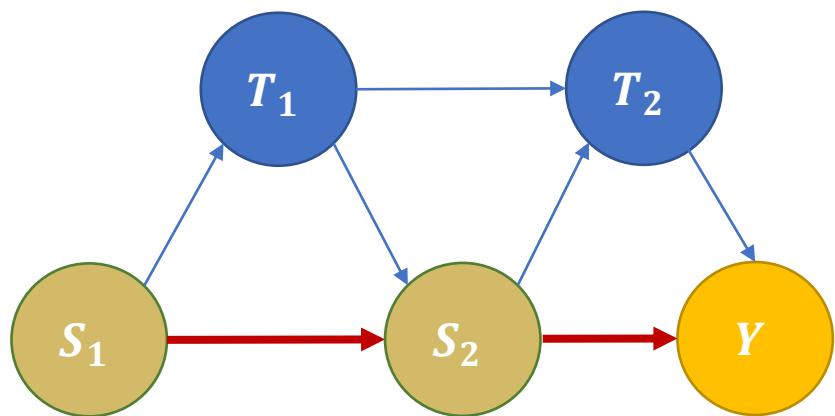
They present a separate set of new problems

- **Correlation:** samples stemming from each unit (or sometimes cluster of units from the same “site”) are correlated
- **Censoring/Missingness:** units typically drop out at (potentially not) random times before the end of the study
- **Dynamic treatments:** units are treated with multiple treatments over time in a manner that is adaptive and auto-correlated

Dynamic Treatment Regime

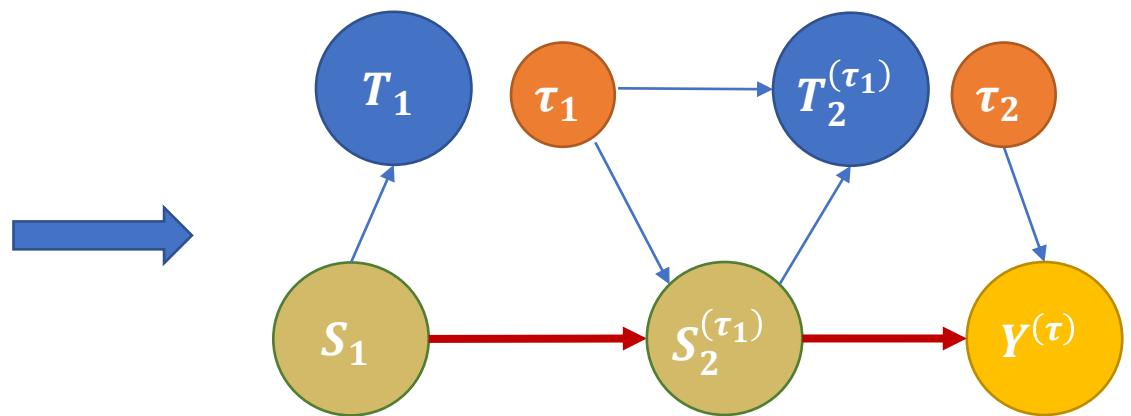
Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

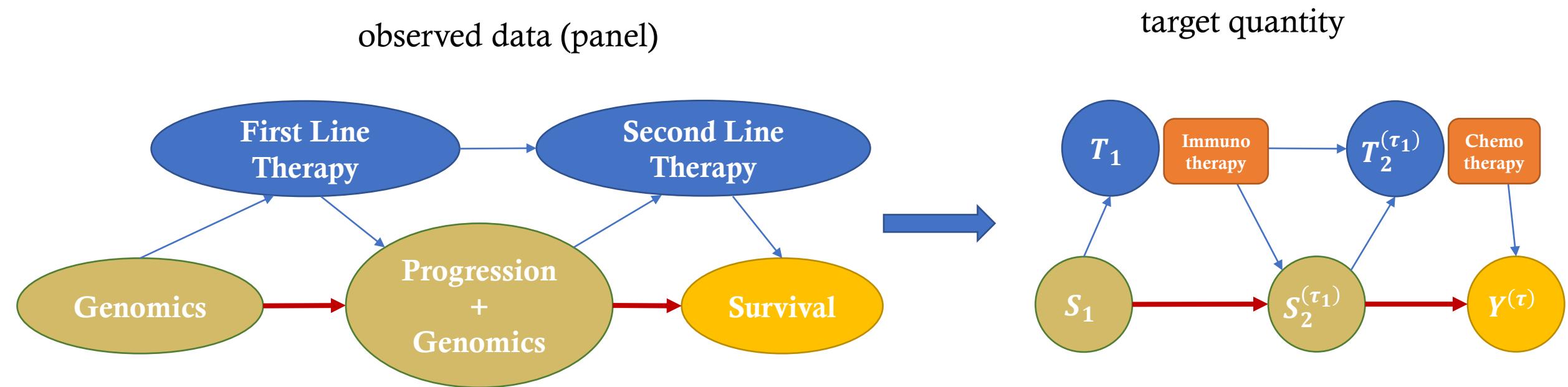
$$\theta := E[Y^{(\tau)}]$$



- ❖ Treatments are offered in an adaptive manner, in response to previous period controls
- ❖ The surrogate – control feedback precludes viewing this as a one-shot treatment problem
- ❖ Setting is known as the dynamic treatment regime [Robins'94,'04, Chakraborty-Murphy'14]

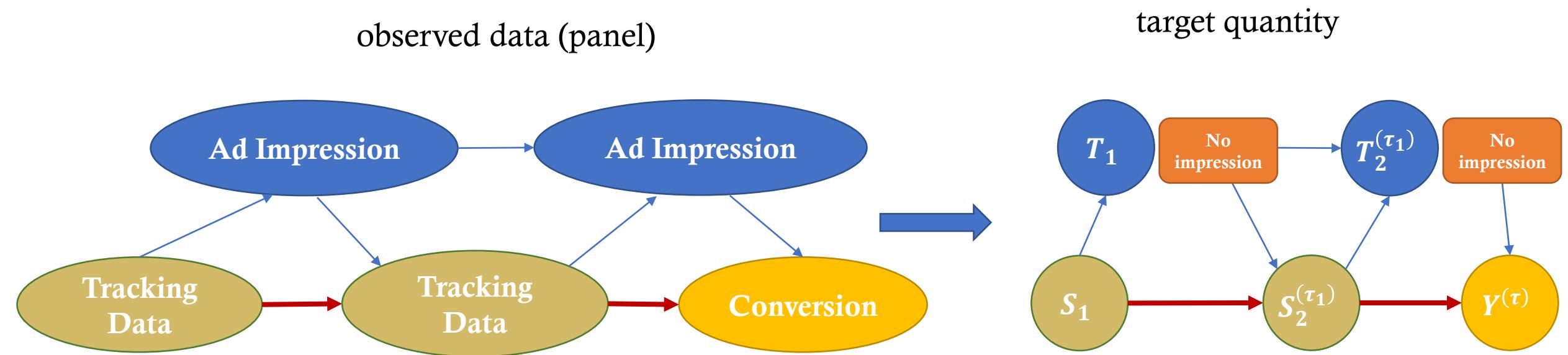
Examples

- Healthcare: Patients treated over time and adaptively



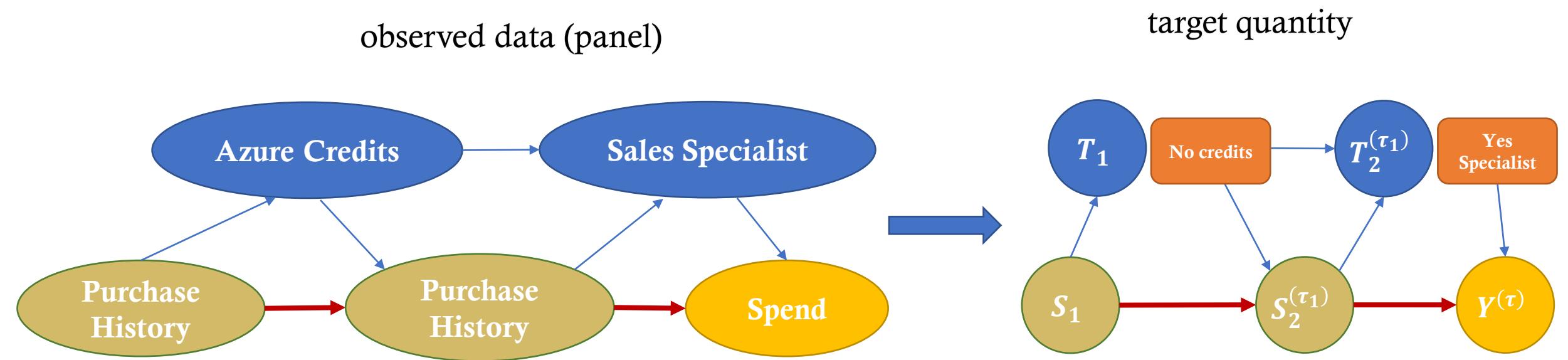
Examples

- Digital Marketing: Web users shown ads multiple times



Examples

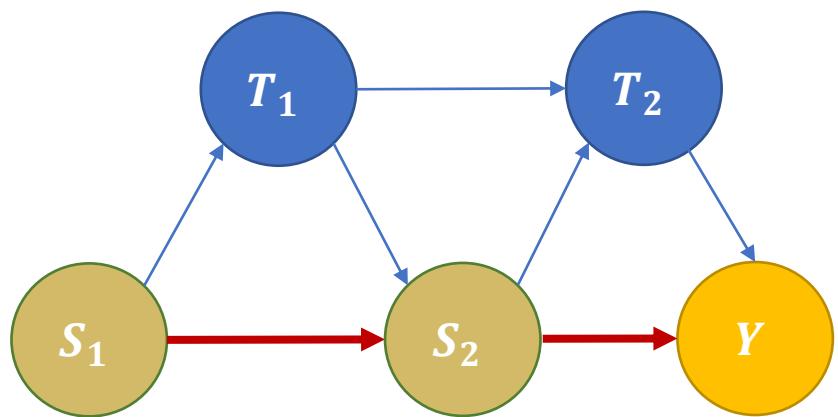
- Business-to-Business (B2B) Operations: Customers offered multiple discount/support interventions over time



Identification in the Dynamic Treatment Regime

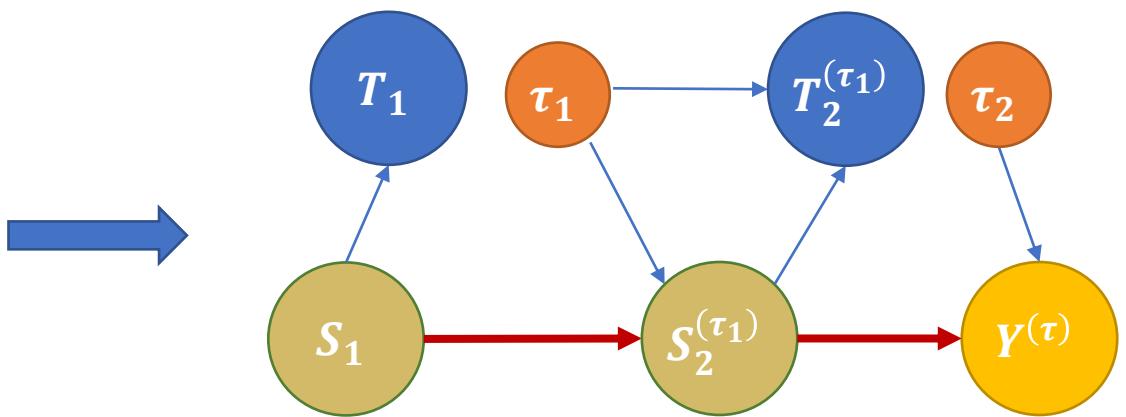
Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

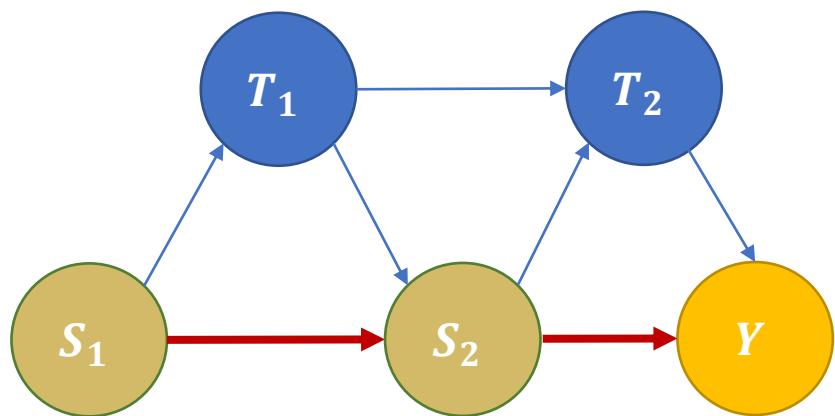
$$\theta := E[Y^{(\tau)}]$$



- ❖ Treatments are offered in an adaptive manner, in response to previous period controls
- ❖ The surrogate – control feedback precludes viewing this as a one-shot treatment problem
- ❖ Setting is known as the dynamic treatment regime [Robins'94,'04, Chakraborty-Murphy'14]

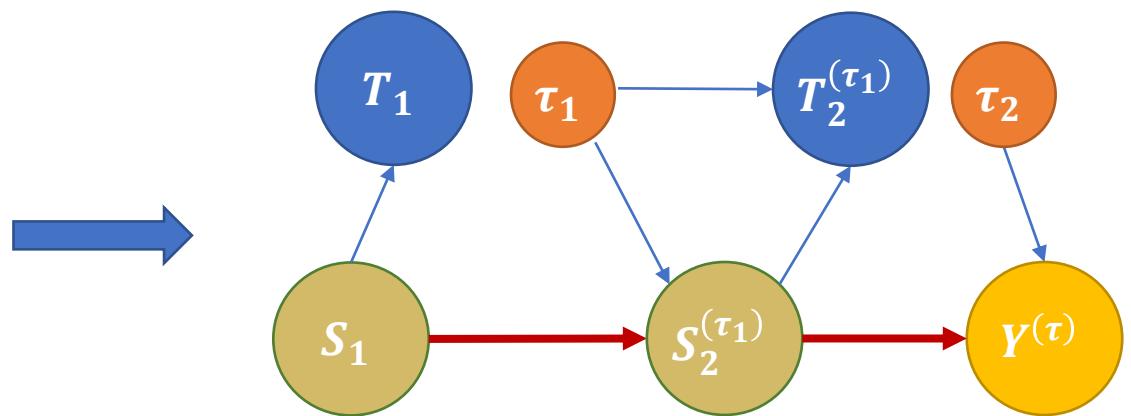
Identification in Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

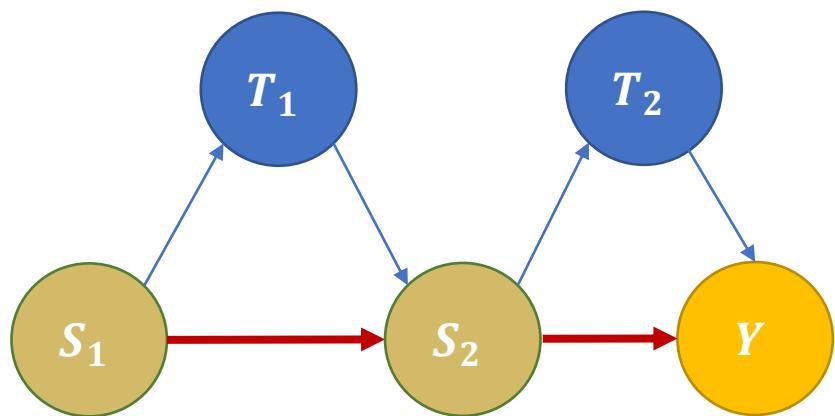
$$\theta := E[Y^{(\tau)}]$$



- ❖ Since there is no unobserved confounding, why not just estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on $S = (S_1, S_2)$
- ❖ Wrong: conditioning on S_2 blocks all the effect from T_1 !

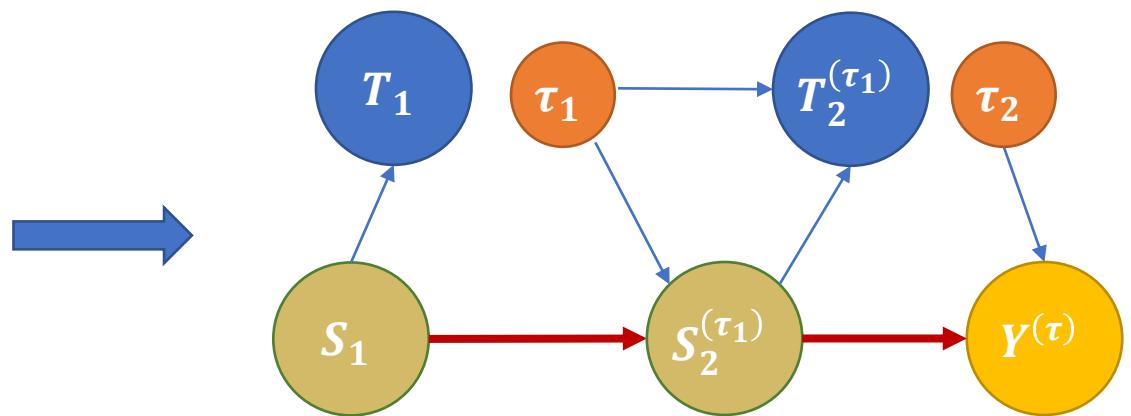
Identification in Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

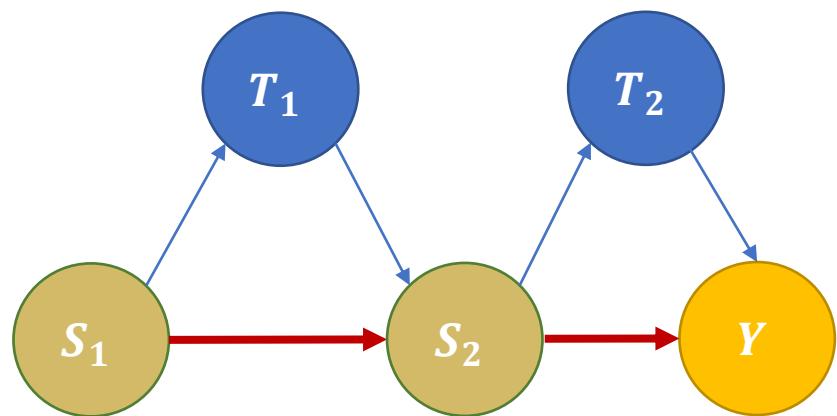
$$\theta := E[Y^{(\tau)}]$$



- ❖ Since there is no unobserved confounding, why not just estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on $S = (S_1, S_2)$
- ❖ Wrong: conditioning on S_2 blocks all the effect from T_1 !

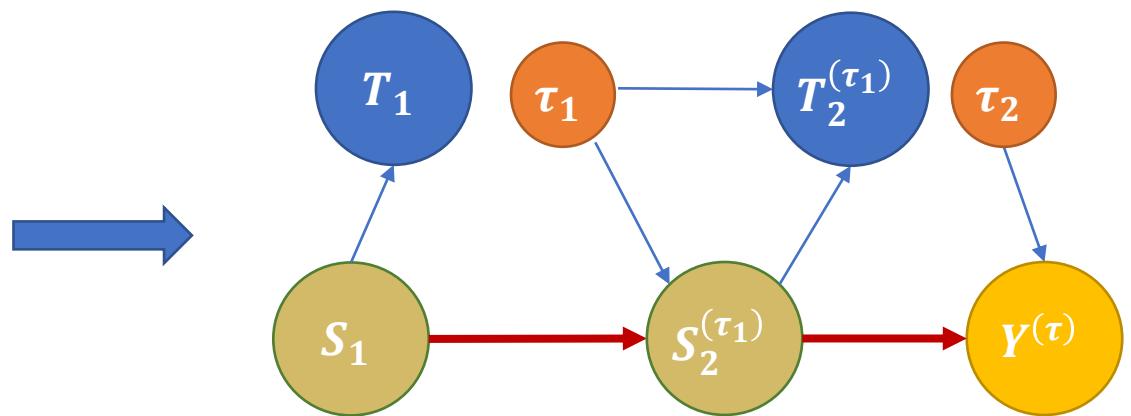
Identification in Dynamic Treatment Regime

observed data (panel)



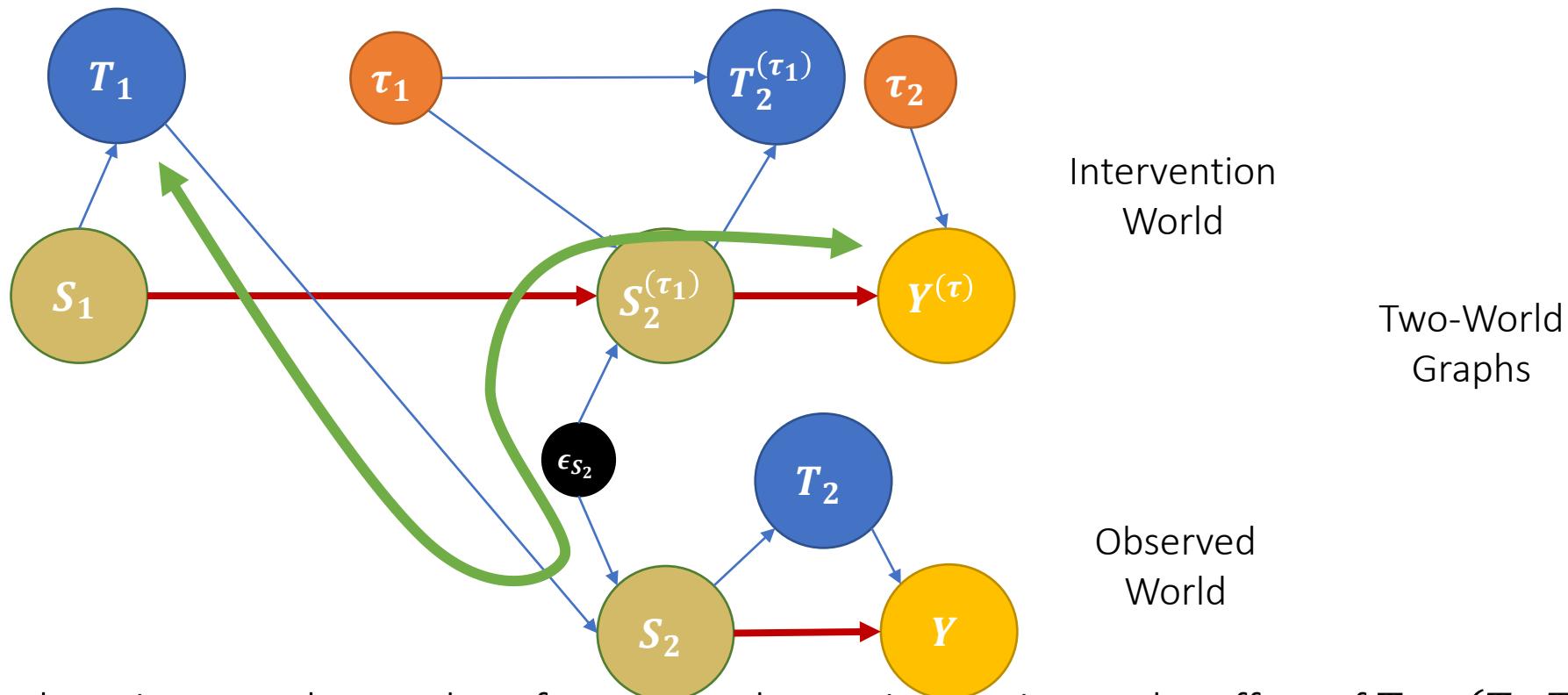
target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

$$\theta := E[Y^{(\tau)}]$$



- ❖ Since there is no unobserved confounding, why not just estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on $S = (S_1, S_2)$
- ❖ Wrong: conditioning on S_2 blocks all the effect from T_1 !
- ❖ Can also be understood using D-separation that $Y^{(\tau)} \perp\!\!\!\perp T_1 \mid S_1, S_2$

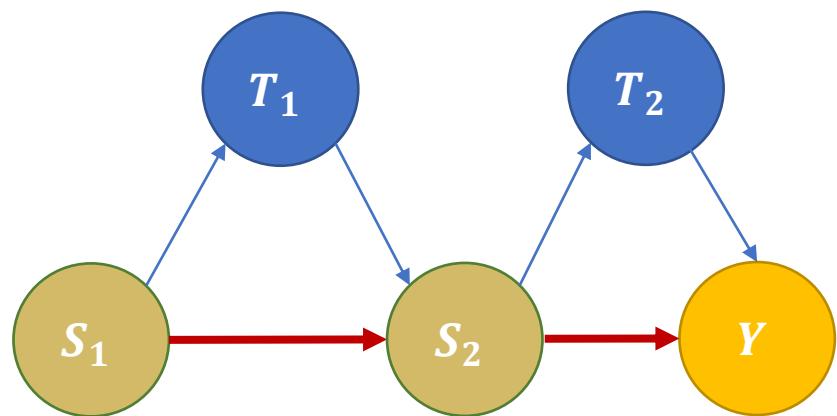
Identification in Dynamic Treatment Regime



- ❖ Since there is no unobserved confounding, why not just estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on $S = (S_1, S_2)$
- ❖ Wrong: conditioning on S_2 blocks all the effect from T_1 !
- ❖ Can also be understood using D-separation that $Y^{(\tau)} \perp\!\!\!\perp T_1 \mid S_1, S_2$

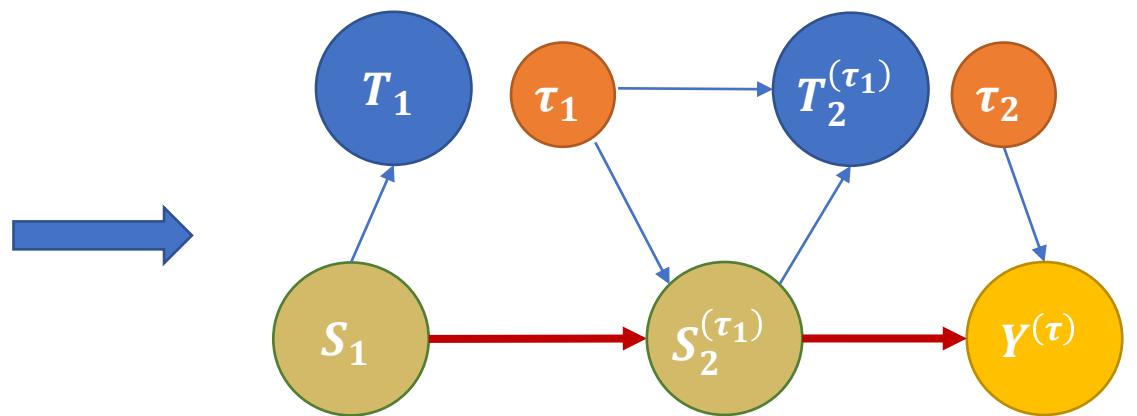
Identification in Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

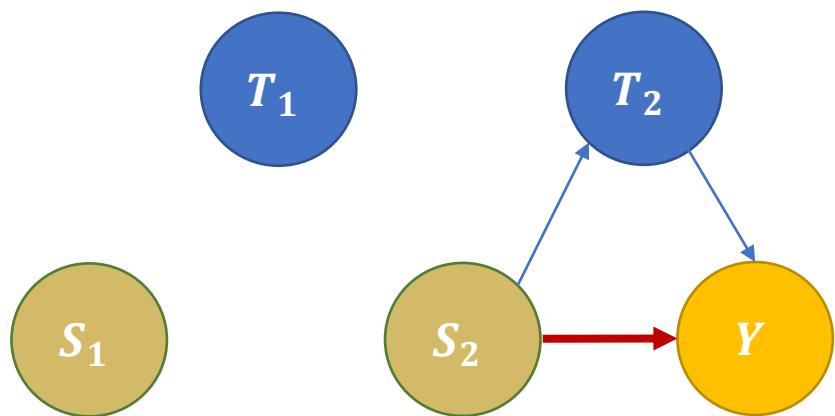
$$\theta := E[Y^{(\tau)}]$$



- ❖ Then let's condition only on S_1 , i.e. estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on S_1
- ❖ Wrong: conditioning only on S_1 leaves “unobserved confounding” between T_2 and Y !

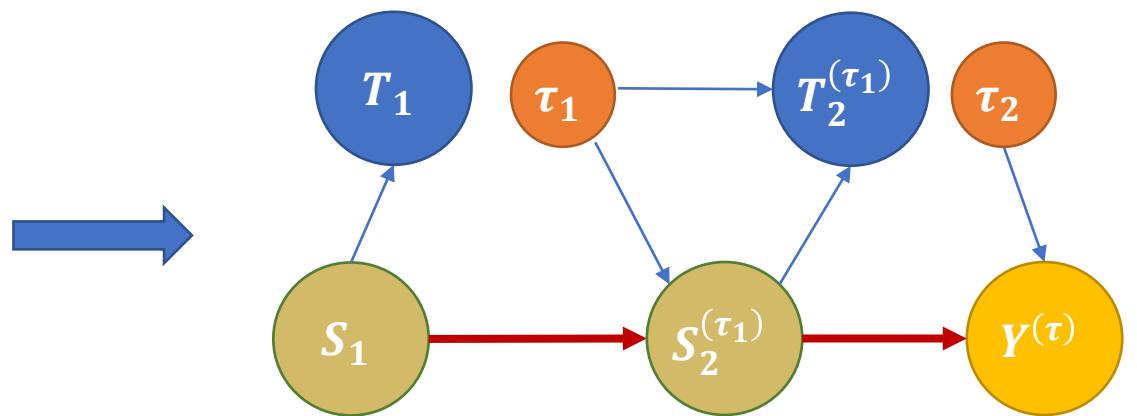
Identification in Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

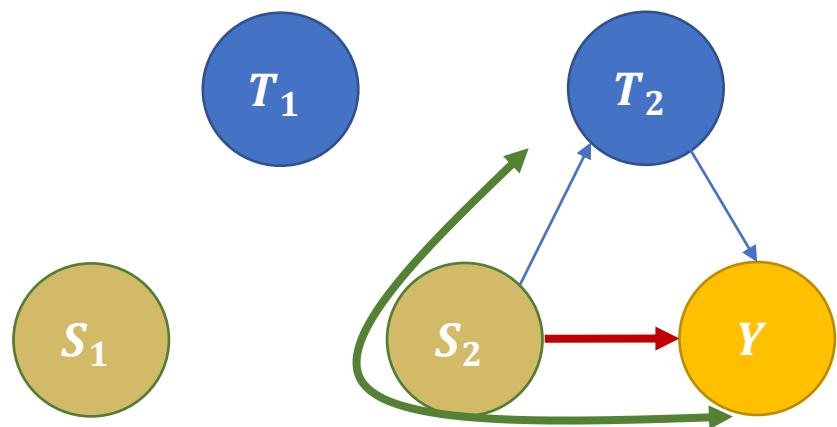
$$\theta := E[Y^{(\tau)}]$$



- ❖ Then let's condition only on S_1 , i.e. estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on S_1
- ❖ Wrong: conditioning only on S_1 leaves “unobserved confounding” between T_2 and Y !

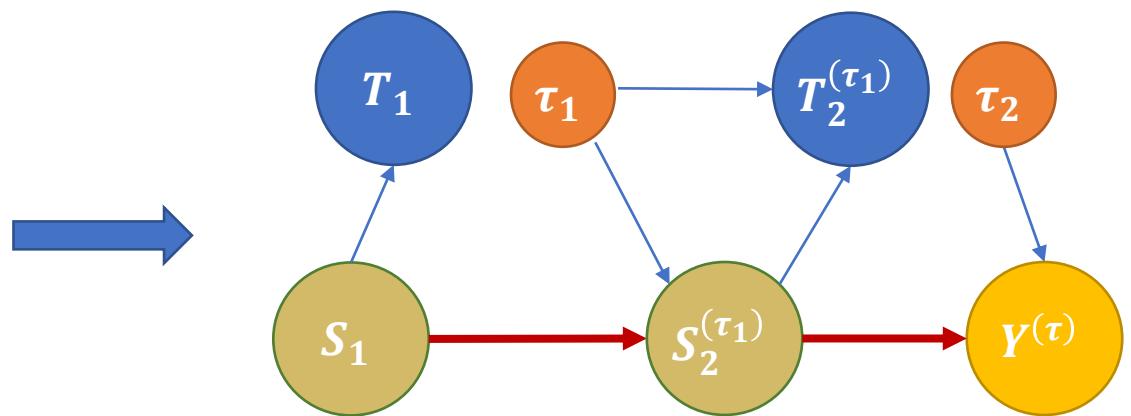
Identification in Dynamic Treatment Regime

observed data (panel)



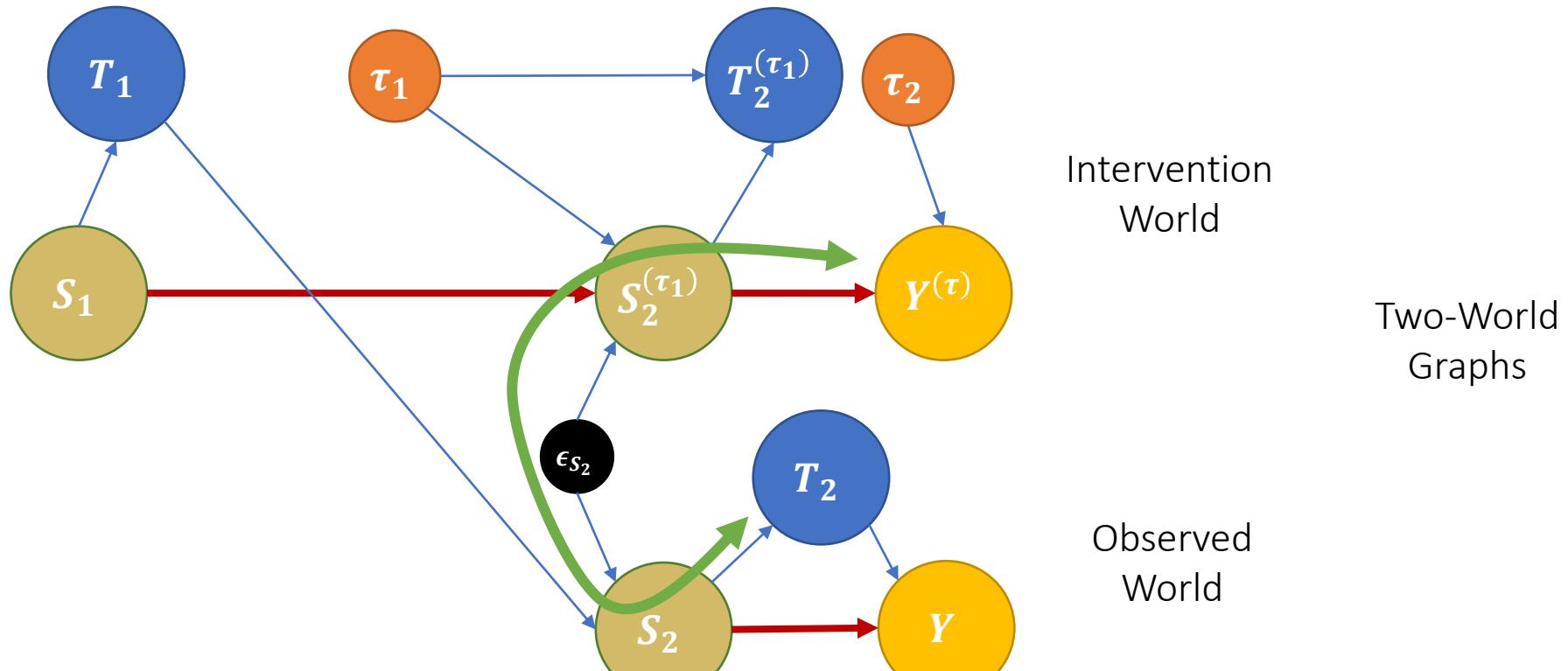
target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

$$\theta := E[Y^{(\tau)}]$$



- ❖ Then let's condition only on S_1 , i.e. estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on S_1
- ❖ Wrong: conditioning only on S_1 leaves “unobserved confounding” between T_2 and Y !

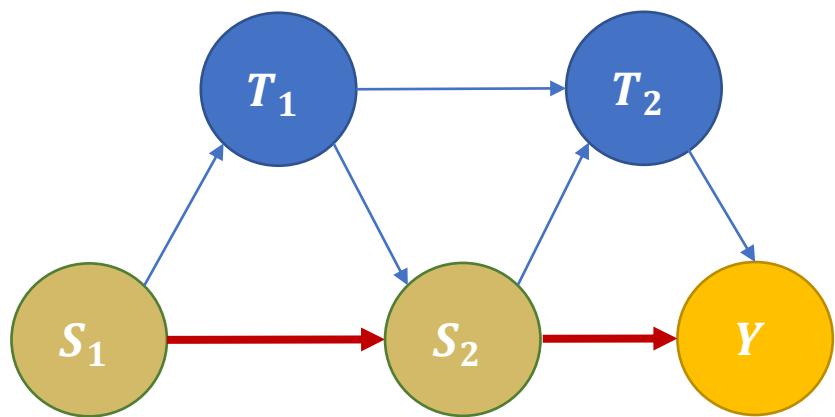
Identification in Dynamic Treatment Regime



- ◇ Then let's condition only on S_1 , i.e. estimate the effect of $T = (T_1, T_2)$ by conditioning, conditioning on S_1
- ◇ Wrong: conditioning only on S_1 leaves “unobserved confounding” between T_2 and Y !
- ◇ Can also be understood using D-separation that $Y^{(\tau)} \perp\!\!\!\perp T_2 | S_1$

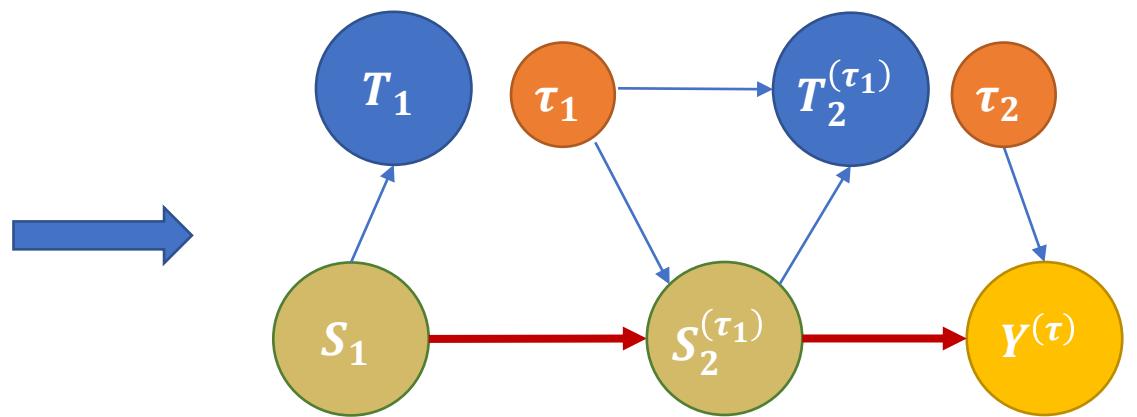
Dynamic Treatment Regime

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

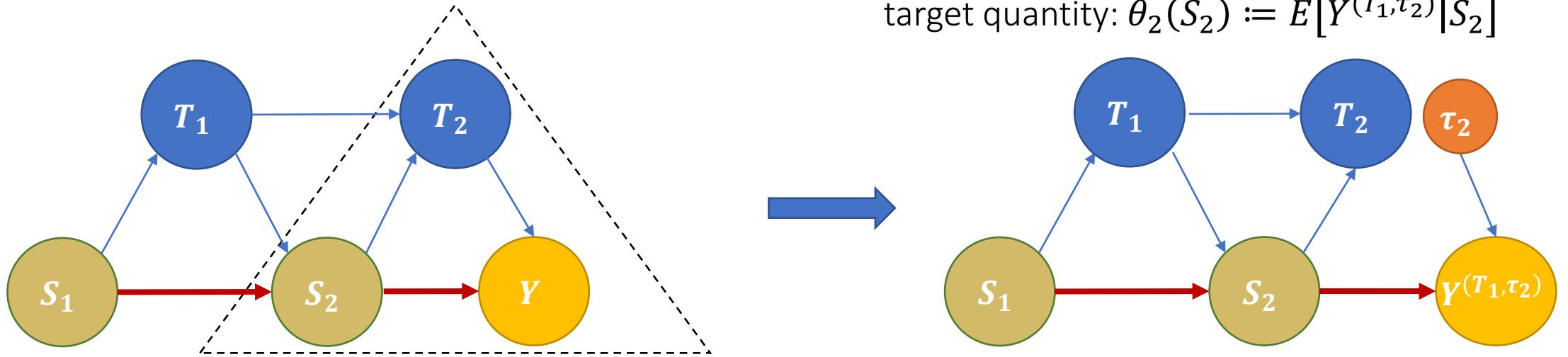
$$\theta := E[Y^{(\tau)}]$$



- ❖ We cannot identify this expected potential outcome simply by conditioning!
- ❖ Let's take it one step at a time
- ❖ What can we identify by conditioning?

Identification of Last Period Intervention

target quantity: $\theta_2(S_2) := E[Y^{(T_1, \tau_2)} | S_2]$



We have conditional ignorability $Y^{(T_1, \tau_2)} \perp\!\!\!\perp T_2 | S_2$

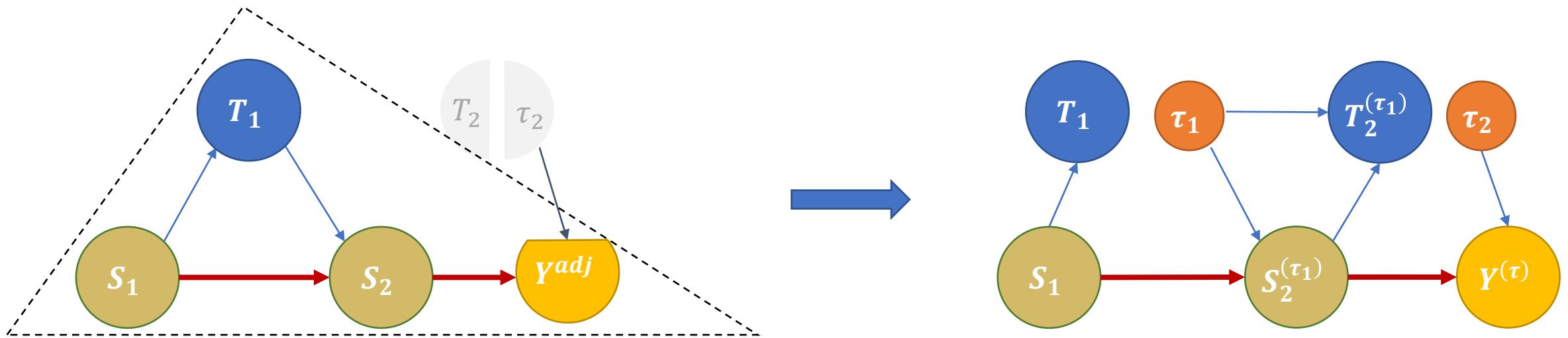
Target quantity $\theta_2(S_2)$ can be identified by conditioning

$$\begin{aligned}\theta_2(S_2) &= E[Y^{(T_1, \tau_2)} | S_2] \\ &= E[Y^{(T_1, \tau_2)} | T_2 = \tau_2, S_2] \\ &= E[Y | T_2 = \tau_2, S_2]\end{aligned}$$

Identification by conditioning:

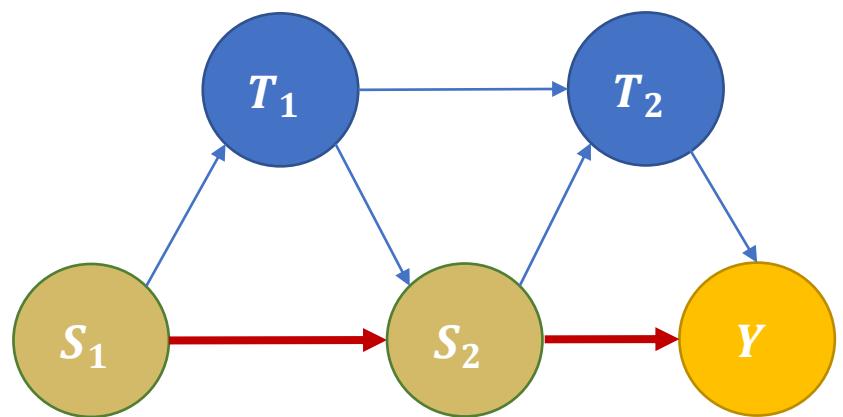
- Train a predictive model $f_2(T_2, S_2)$
 $Y \sim T_2, S_2$
- Evaluate the model at $T_2 = \tau_2$
 $\theta_2(S_2) = f_2(\tau_2, S_2)$

Identification via Backwards Recursion

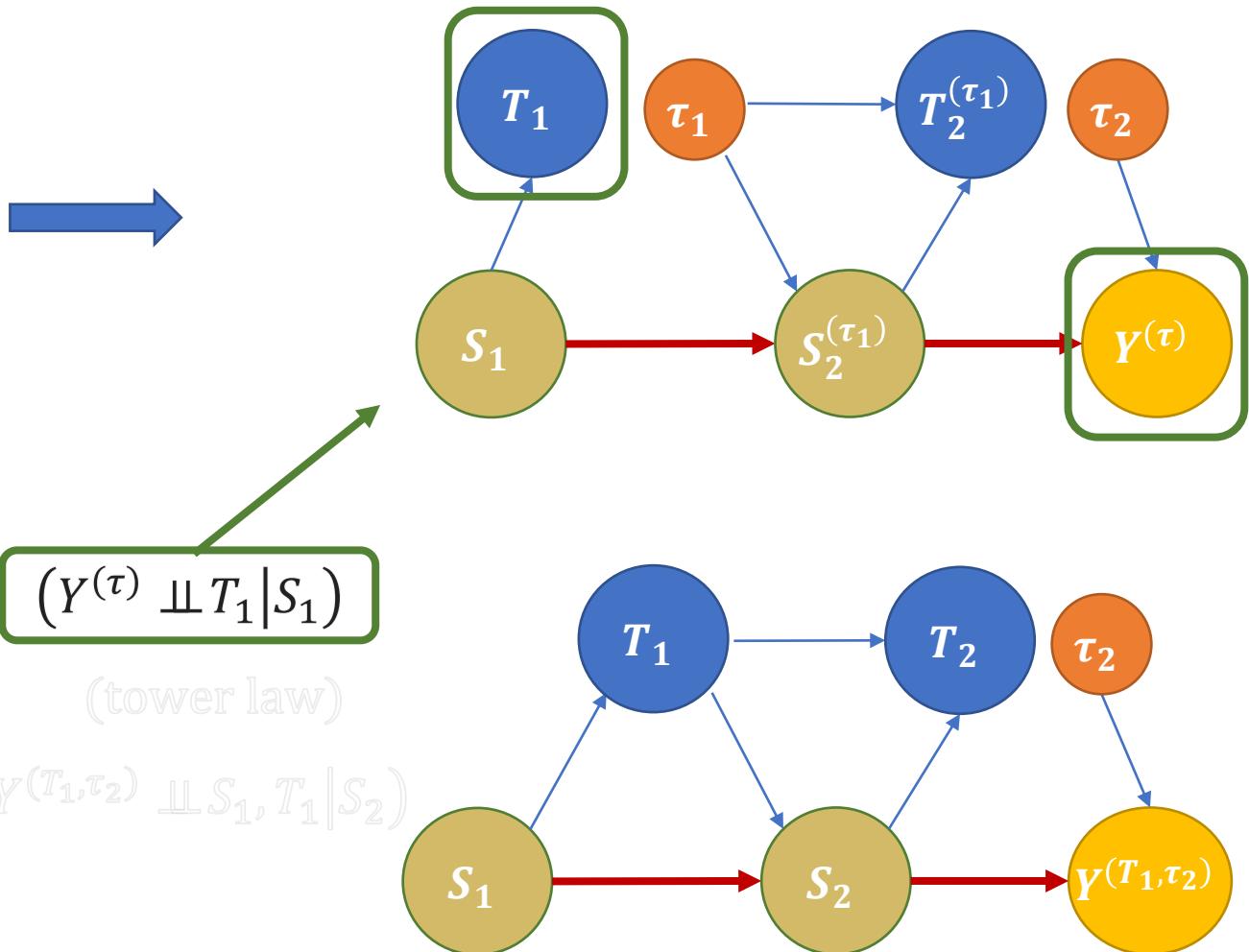


- But now we can “adjust the observed outcome” using this counterfactual model to remove the second period treatment!
- Replace Y with $Y_{adj} := f_2(\tau_2, S_2)$
- Estimating the target quantity $E[Y^{(\tau_1, \tau_2)}]$ is the same as estimating the effect of T_1 on Y_{adj}
- We can estimate this by conditioning on S_1

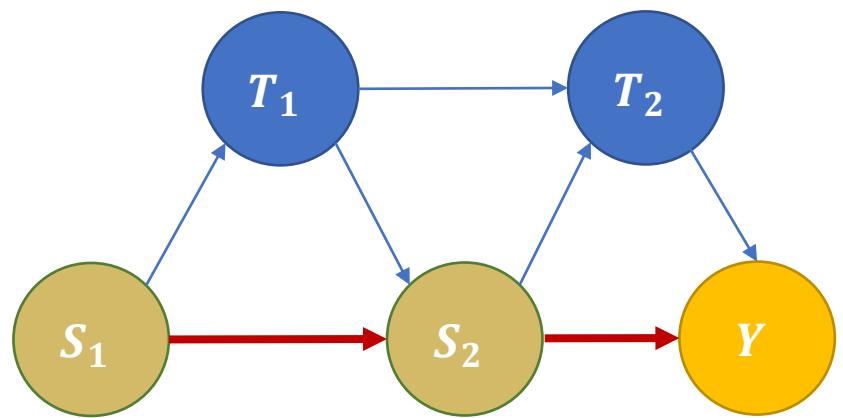
Identification via Backwards Recursion



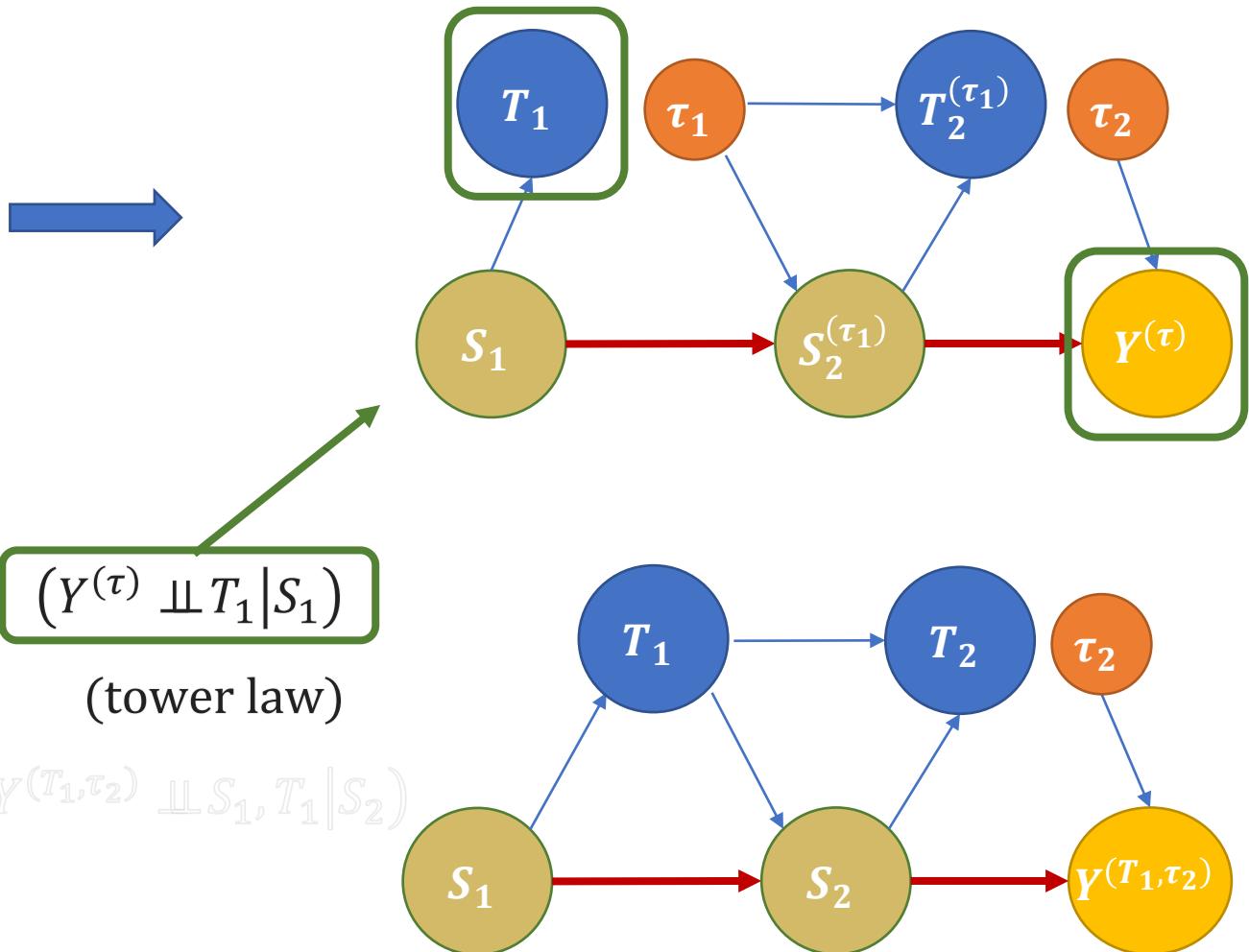
$$\begin{aligned}
 E[Y^{(\tau_1, \tau_2)}] &= E\left[E\left[Y^{(\tau_1, \tau_2)}|S_1\right]\right] \\
 &= E\left[E\left[Y^{(T_1, \tau_2)}|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)}|S_2, T_1, S_1\right]|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)}|S_2\right]|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[\theta_2(S_2)|T_1 = \tau_1, S_1\right]\right]
 \end{aligned}$$



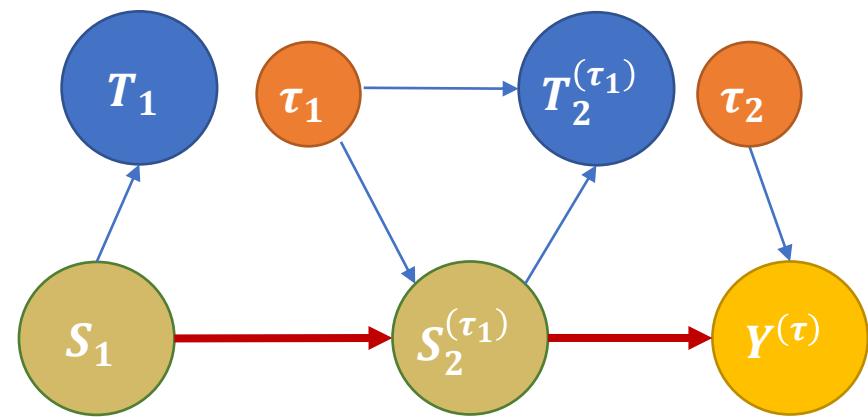
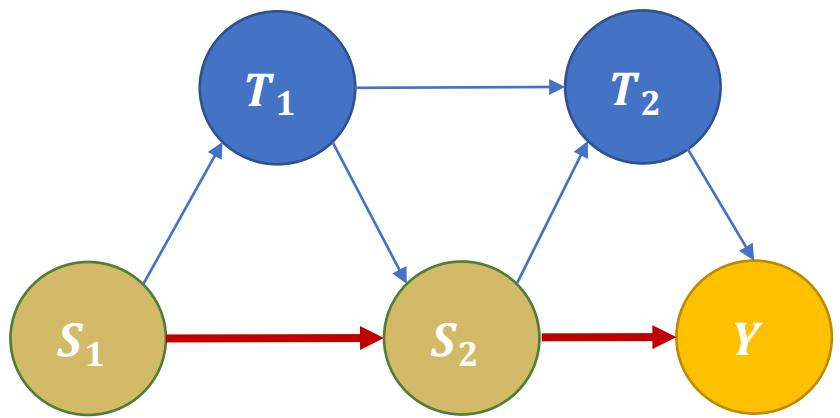
Identification via Backwards Recursion



$$\begin{aligned}
 E[Y^{(\tau_1, \tau_2)}] &= E\left[E\left[Y^{(\tau_1, \tau_2)}|S_1\right]\right] \\
 &= E\left[E\left[Y^{(T_1, \tau_2)}|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)}|S_2, T_1, S_1\right]|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)}|S_2\right]|T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E[\theta_2(S_2)|T_1 = \tau_1, S_1]\right]
 \end{aligned}$$



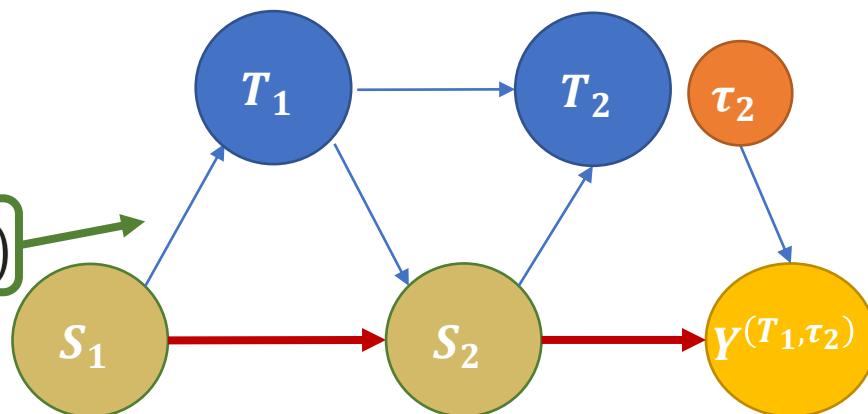
Identification via Backwards Recursion



$$\begin{aligned}
 E[Y^{(\tau_1, \tau_2)}] &= E\left[E\left[Y^{(\tau_1, \tau_2)} | S_1\right]\right] \\
 &= E\left[E\left[Y^{(T_1, \tau_2)} | T_1 = \tau_1, S_1\right]\right] \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)} | S_2, T_1, S_1\right] | T_1 = \tau_1, S_1\right]\right] && (Y^{(\tau)} \perp\!\!\!\perp T_1 | S_1) \\
 &= E\left[E\left[E\left[Y^{(T_1, \tau_2)} | S_2\right] | T_1 = \tau_1, S_1\right]\right] && (\text{tower law}) \\
 &= E\left[E[\theta_2(S_2) | T_1 = \tau_1, S_1]\right]
 \end{aligned}$$

$(Y^{(T_1, \tau_2)} \perp\!\!\!\perp S_1, T_1 | S_2)$

(definition of θ_2)



Identification Process

- Estimate predictive model $f_2: Y \sim T_2, S_2$
- Adjust outcome for second period treatment

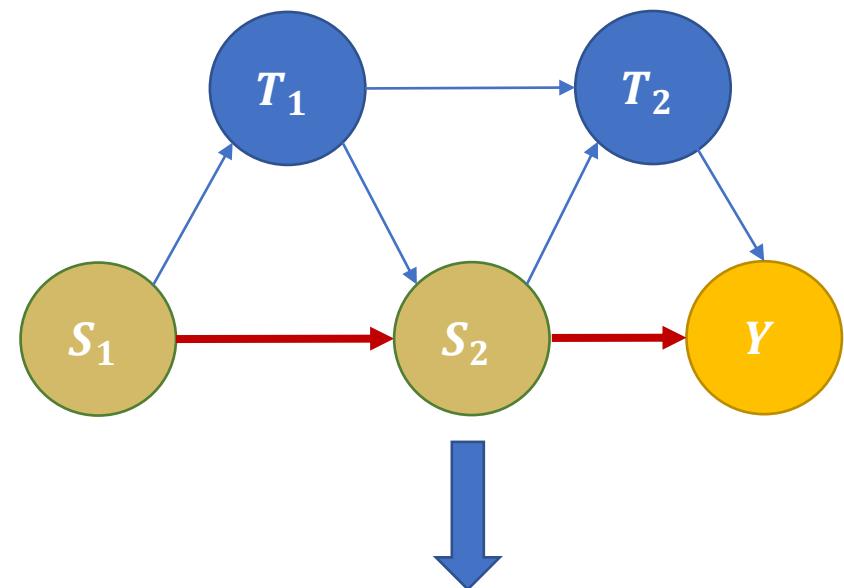
$$Y_{adj} \leftarrow f_2(\tau_2, S_2)$$

- Estimate predictive model $f_1: Y_{adj} \sim T_1, S_1$
- Adjust outcome for first period treatment

$$Y_{adj} \leftarrow f_1(\tau_1, S_1)$$

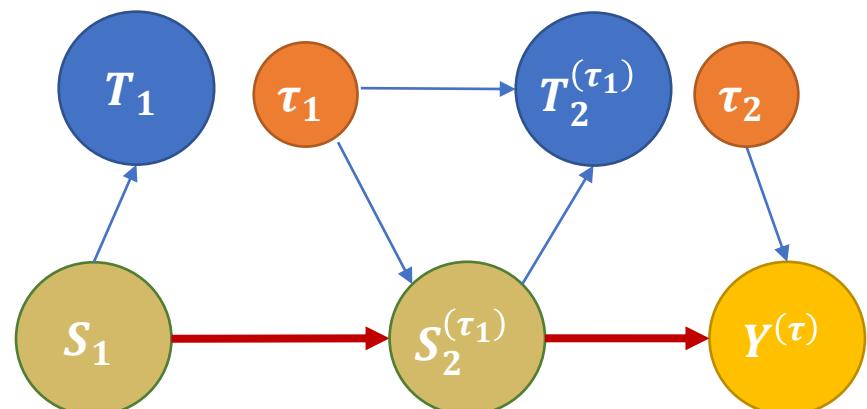
- Return target quantity: $\theta = E[Y_{adj}]$

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

$$\theta := E[Y^{(\tau)}]$$



Moment Based Framework

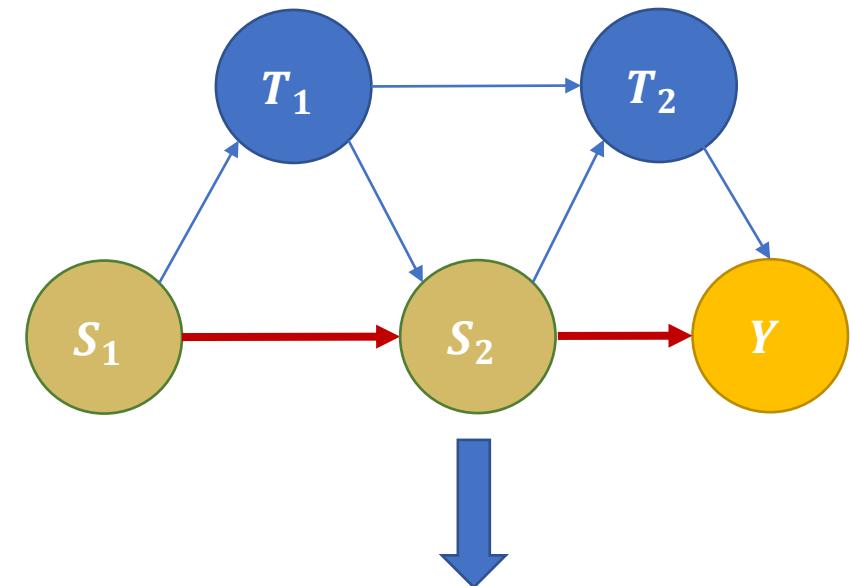
$$\theta = E[f_1(\tau_1, S_1)]$$

$$f_1(T_1, S_1) := E[f_2(\tau_2, S_2) | T_1, S_1]$$

$$f_2(T_2, S_2) := E[Y | T_2, S_2]$$

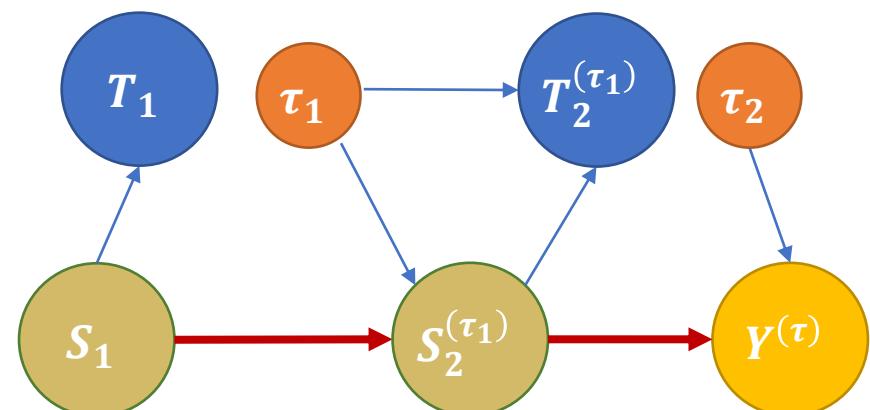
These set of equations (and their generalization to many periods) are known as the *g*-formula

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

$$\theta := E[Y^{(\tau)}]$$



Estimation in the Dynamic Treatment Regime

G-computation

Moment function: $m(Z; \theta, f)$

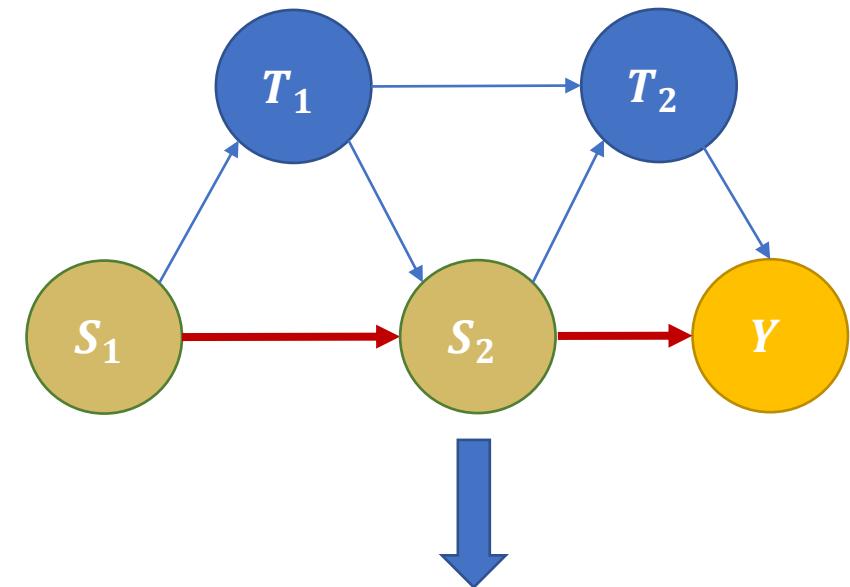
$$E[\theta - f_1(\tau_1, S_1)] = 0$$

$$\begin{aligned} f_1(T_1, S_1) &:= E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &:= E[Y | T_2, S_2] \end{aligned}$$

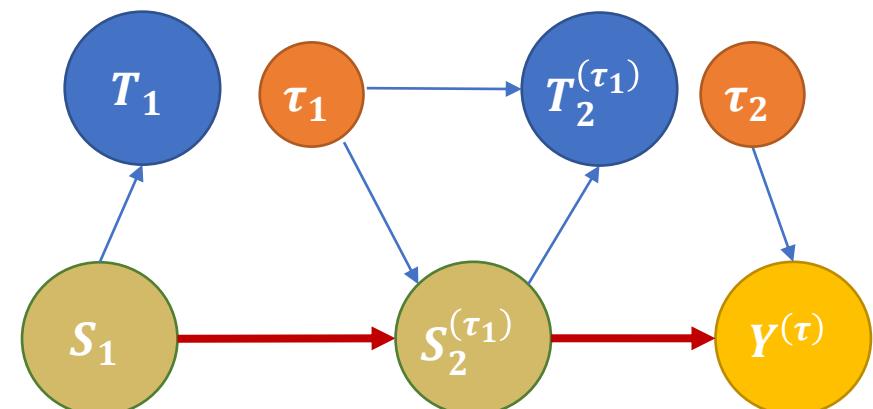
Nuisance functions

- Using parametric models to estimate these functions (e.g. linear or logistic regression) and plug them in the formula is known as the “parametric g computation”

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$
 $\theta := E[Y^{(\tau)}]$



ML Based Estimation

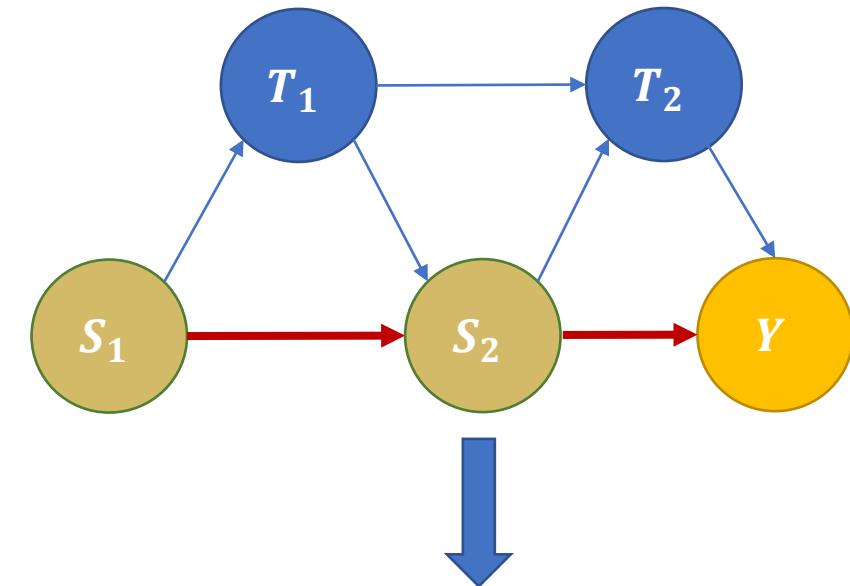
Moment function: $m(Z; \theta, f)$

$$E[\theta - f_1(\tau_1, S_1)] = 0$$

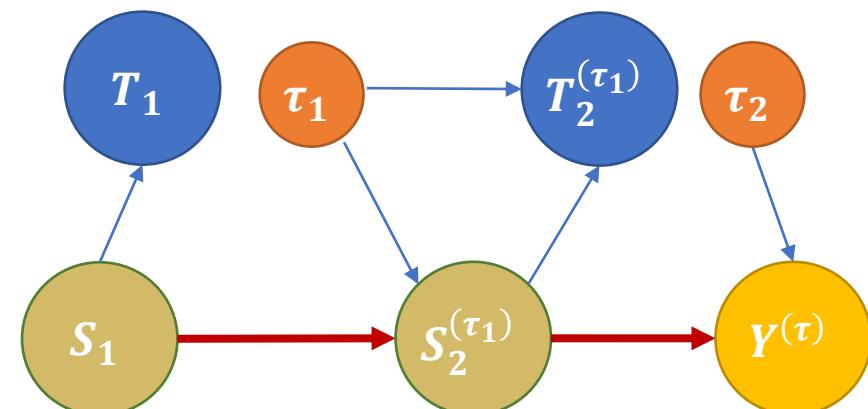
$$\begin{aligned} f_1(T_1, S_1) &:= E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &:= E[Y | T_2, S_2] \end{aligned}$$

- If we use ML for these predictive problems, then we wont be able to construct confidence intervals for θ because moment is not Neyman orthogonal
- We can apply the Inverse Propensity based debiasing idea

observed data (panel)



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$
 $\theta := E[Y^{(\tau)}]$



Debiasing Moments and Neyman Orthogonality

$$E \left[\theta - f_1(\tau_1, S_1) + \frac{1\{\tau_1 = \tau_1\}}{\Pr(\tau_1 = \tau_1 | S_1)} (f_2(\tau_2, S_2) - f_1(\tau_1, S_1)) \right] = 0$$

Original moment

Inverse Propensity

Residual of the Regression Problem
that defines nuisance function f_1

Debiasing correction for
nuisance function f_1

$$\begin{aligned} f_1(T_1, S_1) &:= E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &:= E[Y | T_2, S_2] \end{aligned}$$

Nuisance functions

Debiasing Moments and Neyman Orthogonality

Original moment

$$E \left[\theta - f_1(\tau_1, S_1) + \frac{1\{\tau_1 = \tau_1\}}{\Pr(\tau_1 = \tau_1 | S_1)} (f_2(\tau_2, S_2) - f_1(\tau_1, S_1)) \right] = 0$$

$$\begin{aligned} f_1(T_1, S_1) &:= E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &:= E[Y | T_2, S_2] \end{aligned}$$

Nuisance functions

We are still left with this
nuisance, that we have
not “debiased”

Debiasing correction for
nuisance function f_1

Debiasing Moments and Neyman Orthogonality

We are still left with this nuisance, that we have not “debiased”

Original moment

$$E \left[\theta - f_1(\tau_1, S_1) + \frac{1\{T_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} (f_2(\tau_2, S_2) - f_1(T_1, S_1)) \right] = 0$$

$$+ \frac{1\{T_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} \frac{1\{T_2 = \tau_2\}}{\Pr(T_2 = \tau_2 | S_2)} (Y - f_2(T_2, S_2))$$

$$\begin{aligned} f_1(T_1, S_1) &:= E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &:= E[Y | T_2, S_2] \end{aligned}$$

Debiasing correction for
nuisance function f_2

Nuisance functions

Debiasing Moments and Neyman Orthogonality

Original moment

$$E \left[\theta - f_1(\tau_1, S_1) + \frac{1\{T_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} (f_2(\tau_2, S_2) - f_1(T_1, S_1)) \right] = 0$$

IPS term already
multiplying f_2

$$+ \frac{1\{T_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} \frac{1\{T_2 = \tau_2\}}{\Pr(T_2 = \tau_2 | S_2)} (Y - f_2(T_2, S_2))$$

We are still left with this
nuisance, that we have
not “debiased”

$$\begin{aligned} f_1(T_1, S_1) &\coloneqq E[f_2(\tau_2, S_2) | T_1, S_1] \\ f_2(T_2, S_2) &\coloneqq E[Y | T_2, S_2] \end{aligned}$$

New IPS term for second
period treatment
introduced to transform
 $f_2(T_2, S_2) \rightarrow f_2(\tau_2, S_2)$

Residual of the
regression problem that
defines f_2

Nuisance functions

Debiasing Moments and Neyman Orthogonality

Original moment

$$E \left[\theta - f_1(\tau_1, S_1) + \frac{1\{\tau_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} (f_2(\tau_2, S_2) - f_1(T_1, S_1)) \right] = 0$$

$$+ \frac{1\{\tau_1 = \tau_1\}}{\Pr(T_1 = \tau_1 | S_1)} \frac{1\{\tau_2 = \tau_2\}}{\Pr(T_2 = \tau_2 | S_2)} (Y - f_2(T_2, S_2))$$

$$f_1(T_1, S_1) := E[f_2(\tau_2, S_2) | T_1, S_1]$$

$$f_2(T_2, S_2) := E[Y | T_2, S_2]$$

Nuisance functions

- This moment now satisfies Neyman orthogonality with respect to all the nuisance functions
- We need to also estimate the propensity functions via generic ML classification

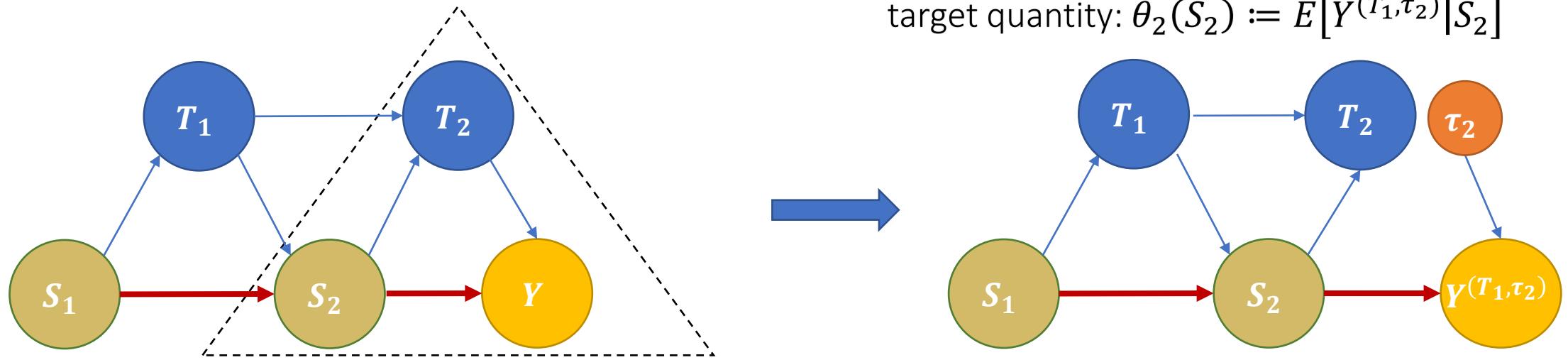
Continuous Treatments and Alternative Approach to Identification and Estimation

What if we have continuous treatments

- The propensity based approach to de-biasing does not apply
- What is the analogue of the “Residual-on-Residual” or “Partialling-Out” (FWL) approach for the dynamic treatment regime, which is also applicable to continuous treatments?

Identification via “Instantaneous” or “Blip” Effects

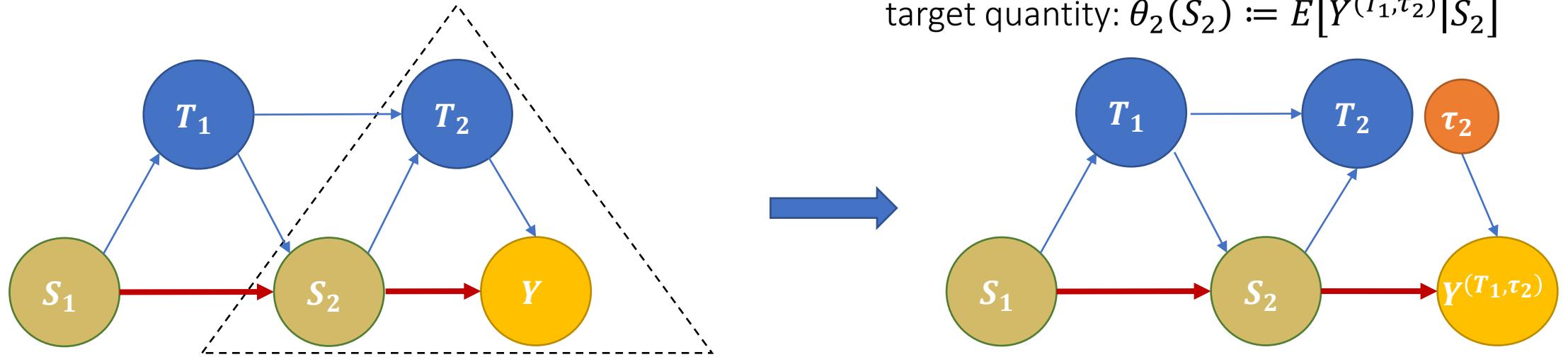
target quantity: $\theta_2(S_2) := E[Y^{(T_1, \tau_2)} | S_2]$



- Instead of estimating the outcome $Y^{(T_1, \tau_2)}$, we will estimate the “effect” of T_2
 $\alpha_2(\tau_2, S_2) := E[Y^{(T_1, \tau_2)} - Y^{(T_1, 0)} | S_2]$
- Adjust outcome by subtracting effect of observed treatment and adding effect of target treatment
 $Y_{\text{adj}} := Y - \alpha_2(T_2, S_2) + \alpha_2(\tau_2, S_2)$

Identification via “Instantaneous” or “Blip” Effects

target quantity: $\theta_2(S_2) := E[Y^{(T_1, \tau_2)} | S_2]$



- Instead of estimating the outcome $Y^{(T_1, \tau_2)}$, we will estimate the “effect” of T_2
$$\alpha_2(\tau_2, S_2) := E[Y^{(T_1, \tau_2)} - Y^{(T_1, 0)} | S_2]$$
- Adjust outcome by subtracting effect of observed treatment and adding effect of target treatment
$$Y_{\text{adj}} := Y - \alpha_2(T_2, S_2) + \alpha_2(\tau_2, S_2)$$
- If “effect” is assumed to have a simple parametric form, then this approach leverages this simplicity!

Repeat in a backwards manner

Identification Process

Estimate the “effect” $\alpha_2(\tau_2, S_2)$

$$\alpha_2(\tau_2, S_2) := E[Y^{(T_1, \tau_2)} - Y^{(T_1, 0)} | S_2]$$

Adjust outcome for second period treatment

$$Y_{\text{adj}} \leftarrow Y - \alpha_2(T_2, S_2) + \alpha_2(\tau_2, S_2)$$

Estimate the effect $\alpha_1(\tau_1, S_1)$

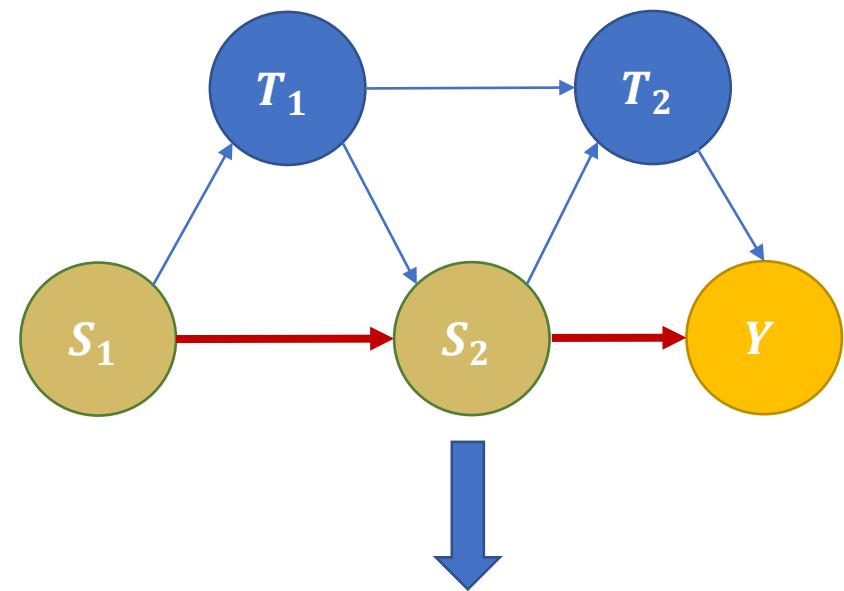
$$E[Y^{(\tau_1, \tau_2)} - Y^{(0, \tau_2)} | S_1] \approx E[Y_{\text{adj}}^{(\tau_1)} - Y_{\text{adj}}^{(0)} | S_1]$$

Adjust outcome for first period treatment

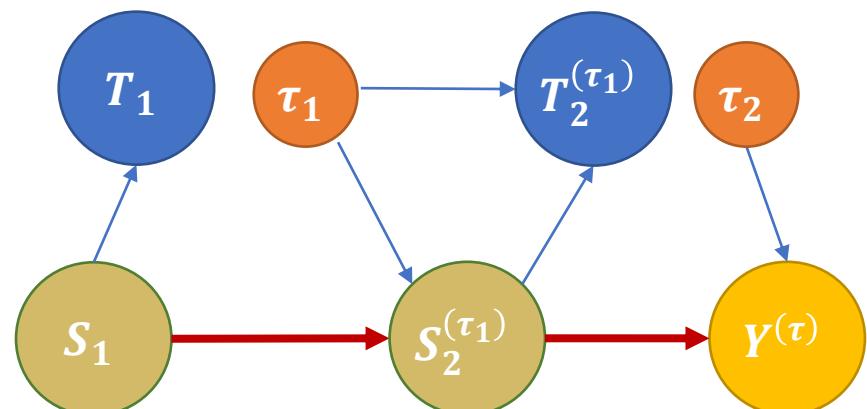
$$Y_{\text{adj}} \leftarrow Y_{\text{adj}} - \alpha_1(T_1, S_1) + \alpha_1(\tau_1, S_1)$$

Return target quantity: $\theta = E[Y_{\text{adj}}]$

observed data (panel)



target quantity: average outcome under a static
treatment sequence (regime) $\tau = (\tau_1, \tau_2)$
 $\theta := E[Y^{(\tau)}]$

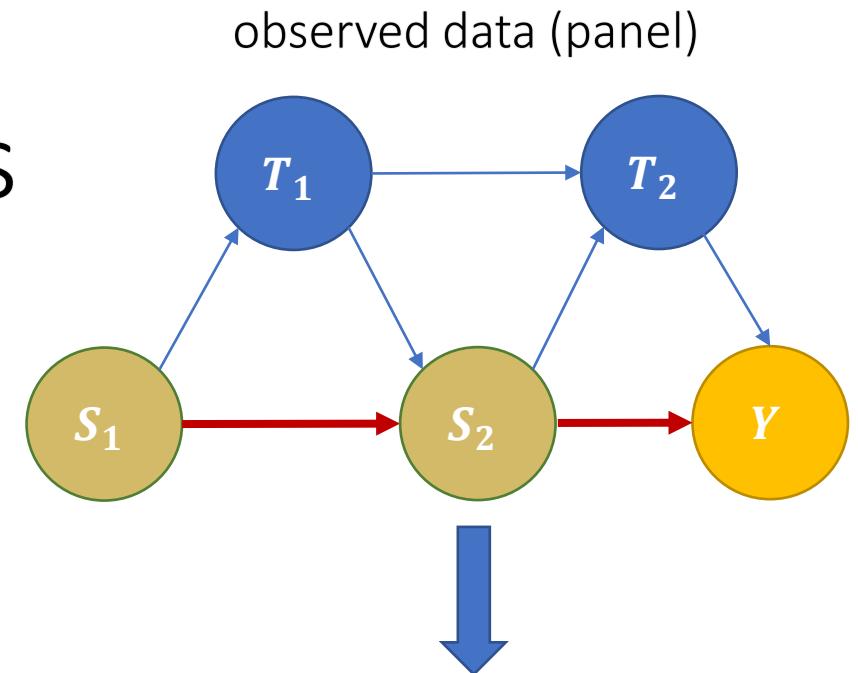


Partial Linearity: Linear Effects

- If we assume that these effects are partially linear:

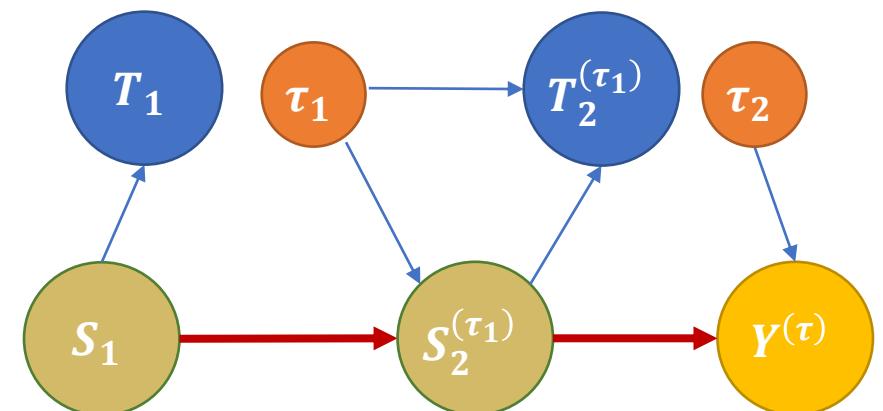
$$a_2(\tau_2, S_2) = \delta_2 \tau_2$$

$$a_1(\tau_1, S_1) := \delta_1 \tau_1$$



target quantity: average outcome under a static treatment sequence (regime) $\tau = (\tau_1, \tau_2)$

$$\theta := E[Y^{(\tau)}]$$



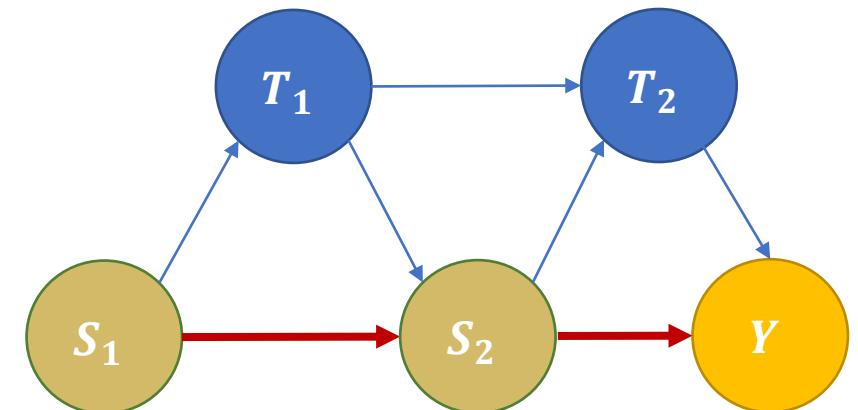
Identification Process

- If we assume that these effects are partially linear:

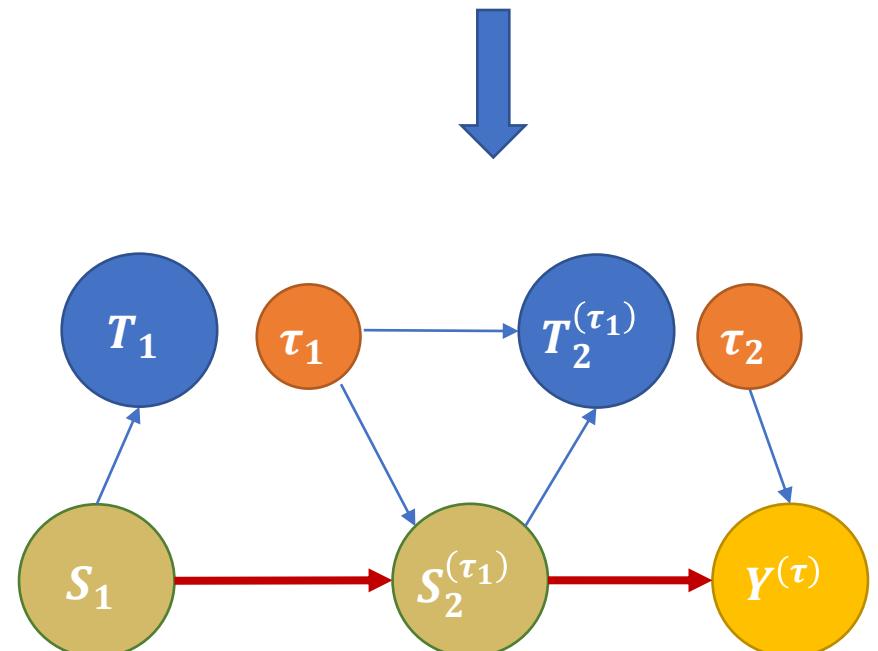
$$a_2(\tau_2, S_2) = \delta_2 \tau_2$$

$$a_1(\tau_1, S_1) := \delta_1 \tau_1$$

observed data (panel)



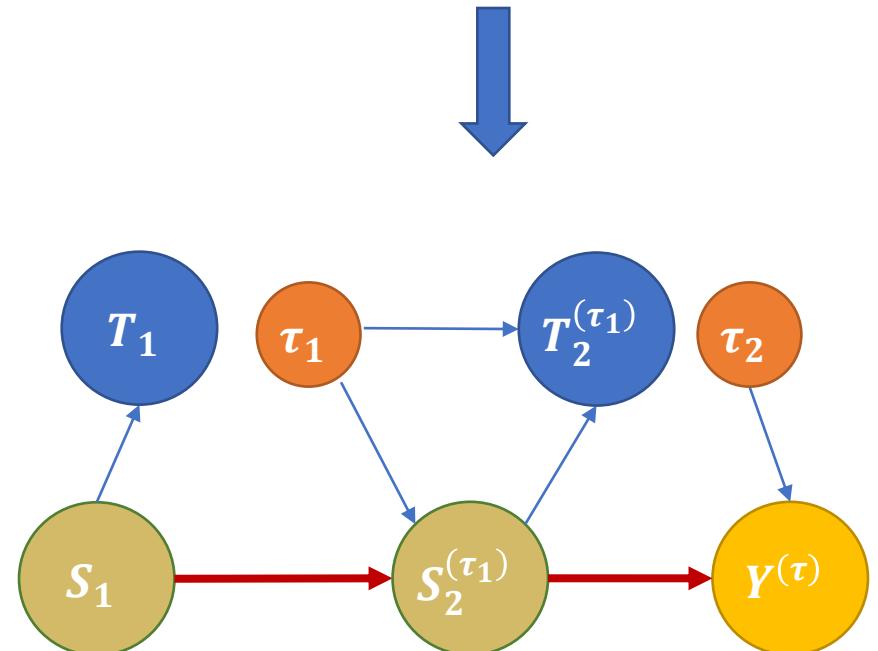
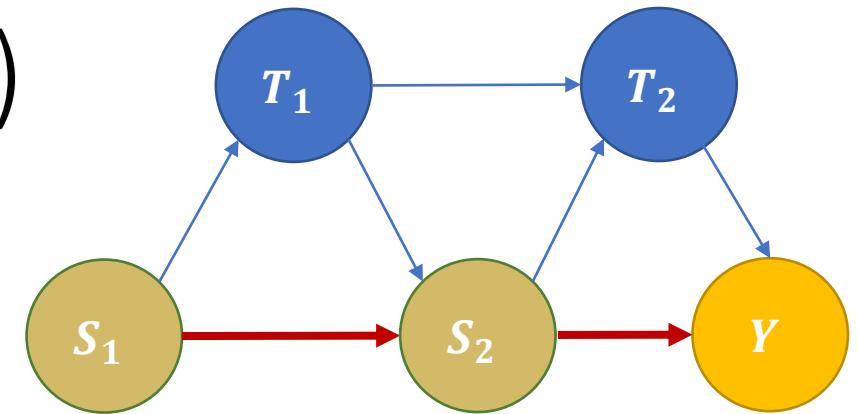
- We can estimate these effects via partialling out!



Dynamic DML (Partialling Out)

- Construct residuals
 $\tilde{Y}_2 := Y - E[Y|S_2], \quad \tilde{T}_2 = T_2 - E[T_2 | S_2]$
- Run OLS: $\tilde{Y}_2 \sim \tilde{T}_2$ to estimate δ_2
- Adjust outcome $Y_{adj} = Y - \delta_2 T_2 + \delta_2 \tau_2$
- Construct residuals:
 $\tilde{Y}_1 := Y_{adj} - E[Y_{adj}|S_1], \quad \tilde{T}_1 := T_1 - E[T_1|S_1]$
- Run OLS: $\tilde{Y}_1 \sim \tilde{T}_1$ to estimate δ_1
- Adjust outcome $Y_{adj} \leftarrow Y_{adj} - \delta_1 T_1 + \delta_1 \tau_1$
- Return: $\theta := E[Y_{adj}]$

observed data (panel)



Moment Based Framework

$$\theta := E[Y - \delta_2 T_2 + \delta_2 \tau_2 - \delta_1 T_1 + \delta_1 \tau_1] = 0$$

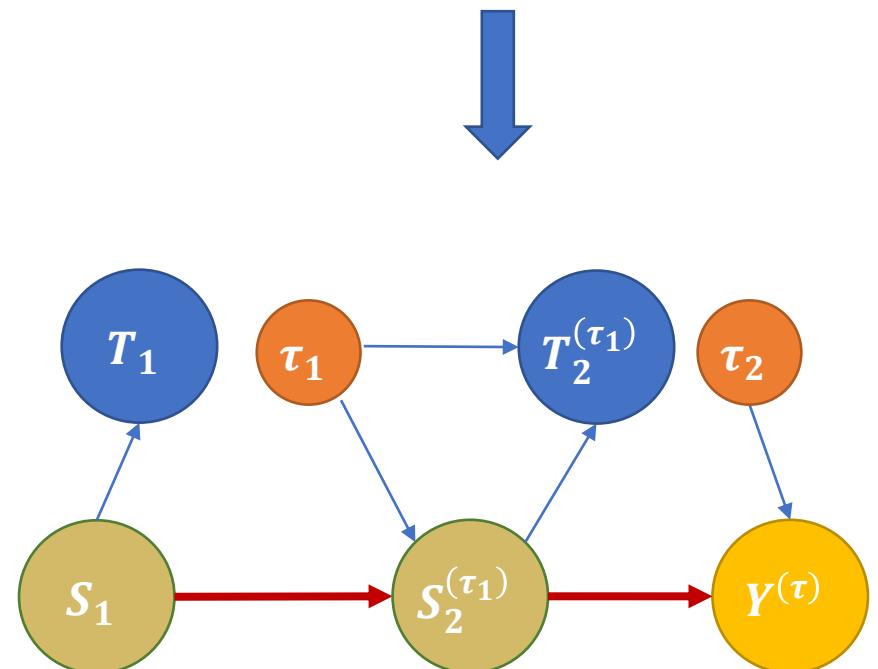
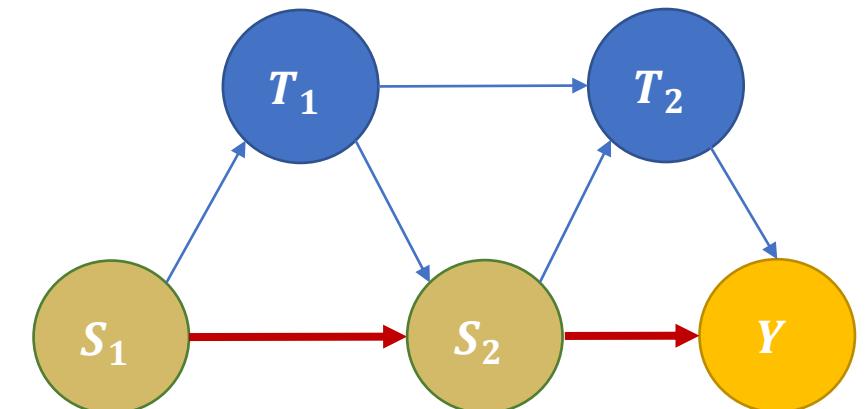
$$E[(\tilde{Y}_2 - \delta_2 \tilde{T}_2) \tilde{T}_2] = 0$$

$$E[(\tilde{Y}_1 - \delta_1 \tilde{T}_1) \tilde{T}_1] = 0$$

$$\begin{aligned} & E[Y|S_2], \quad E[T_2|S_2] \\ & E[Y - \delta_2 T_2 + \delta_2 \tau_2 | S_1], \quad E[T_1|S_1] \end{aligned}$$

Nuisance functions

observed data (panel)



Moment Based Framework

$$\theta := E[Y - \delta_2 T_2 + \delta_2 \tau_2 - \delta_1 T_1 + \delta_1 \tau_1] = 0$$

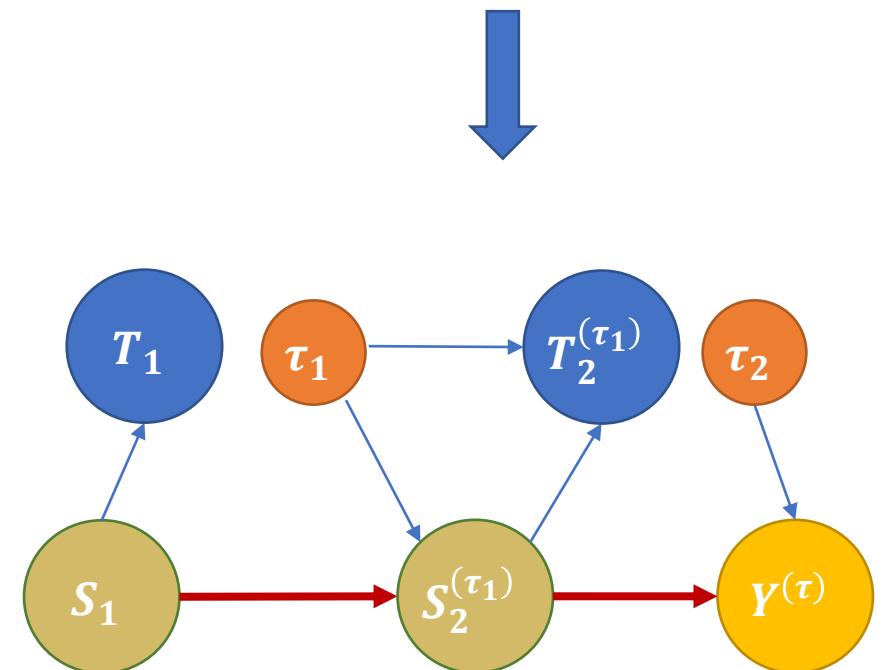
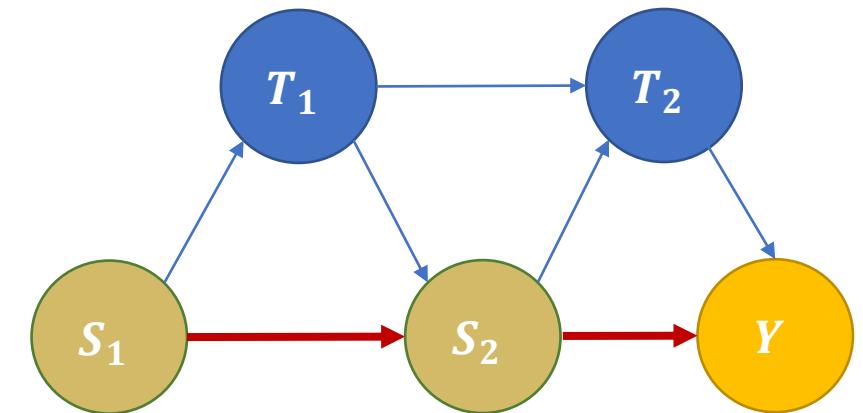
$$E[(\tilde{Y}_2 - \delta_2 \tilde{T}_2) \tilde{T}_2] = 0$$

$$E[(\tilde{Y}_1 - \delta_1 \tilde{T}_1) \tilde{T}_1] = 0$$

$E[Y|S_2], \quad E[T_2|S_2]$
 $E[Y - \delta_2 T_2 + \delta_2 \tau_2|S_1], \quad E[T_1|S_1]$

For the same reason why the Residual-on-Residual moment was Neyman orthogonal, this estimation process is also orthogonal and we can use ML for these problems

observed data (panel)



Long-Term Causal Inference from Short-Term Data with Surrogates

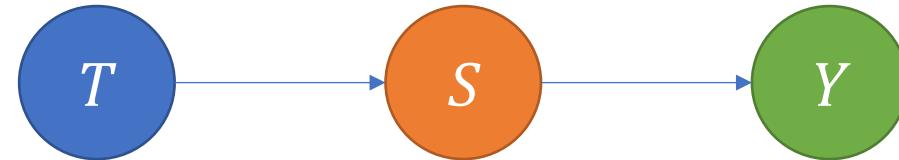
Estimating Long-Term Returns on Investment

- Companies frequently deploys new discount or customer support programs
- Which of these programs (“investments”) are more successful than others?
- Success is a **long-term** objective: what is the effect of the program on the two-year customer journey (e.g., effect on two-year revenue)
- We cannot wait two years to evaluate a program
- **Main Question.** Can we construct estimates of the values of these programs with **short-term** data, e.g. after 6 months?



Long-Term Effects from Short-Term Surrogates

- Suppose that there are many short-term signals S that are indicative of a customer's long-term reward Y (e.g. the next 6-month purchase patterns of a customer could be indicative of their long-term spend)
- Suppose that investment program T affects long-term rewards if and only if it affects these short-term signals



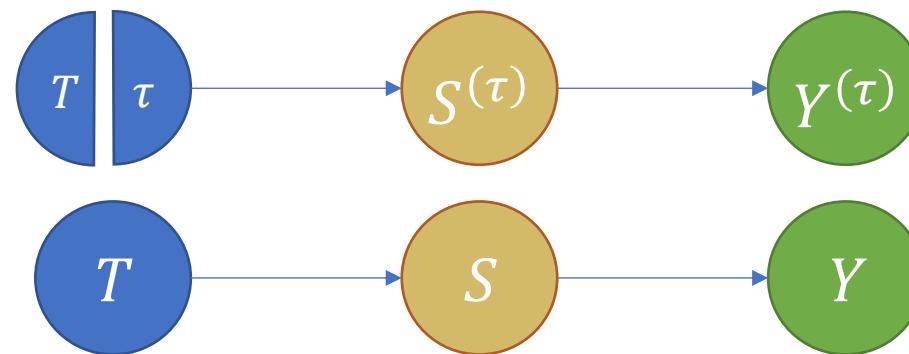
- We will call these short-term signals S surrogates

Causal Inference with Surrogates 101

- Since long-term effect goes only through surrogates:
expected effect on long-term reward = effect on projected long-term reward based on surrogates

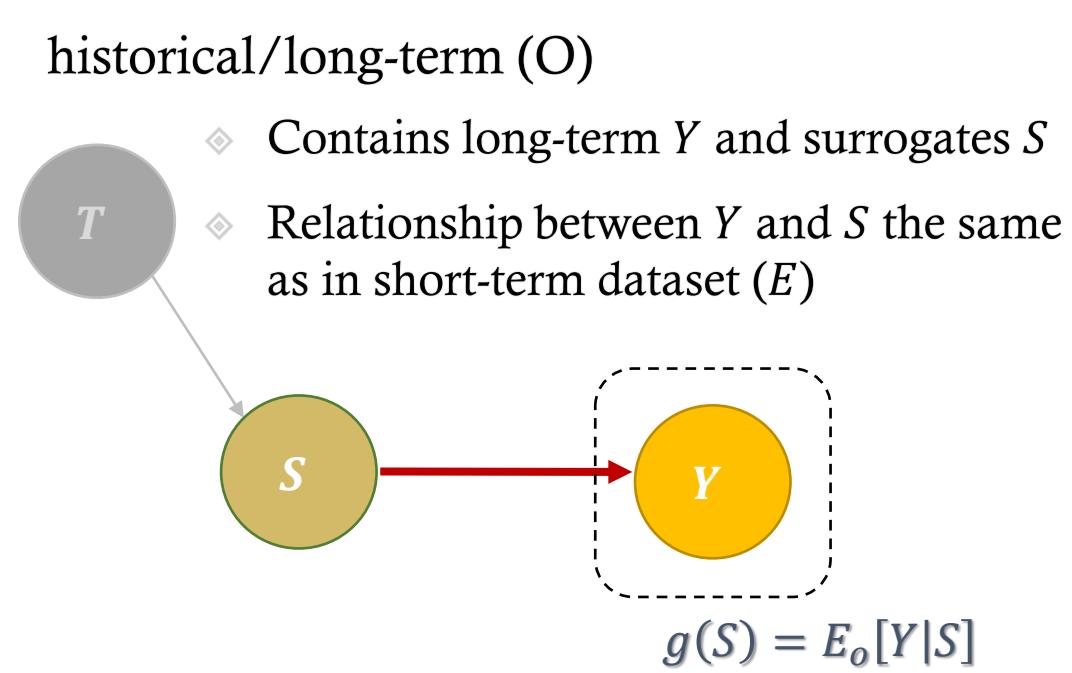
$$E[Y^{(\tau)}] = E[Y^{(\tau)}|T = \tau] = E[Y|T = \tau] = E[E[Y|T = \tau, S]|T = \tau] = E[E[Y|S]|T = \tau]$$

Average reward if intervene and set investment= τ Independence in counterfactual graph Average reward of samples that received investment= τ in data Tower Law of Expectations Forecasted reward from surrogates



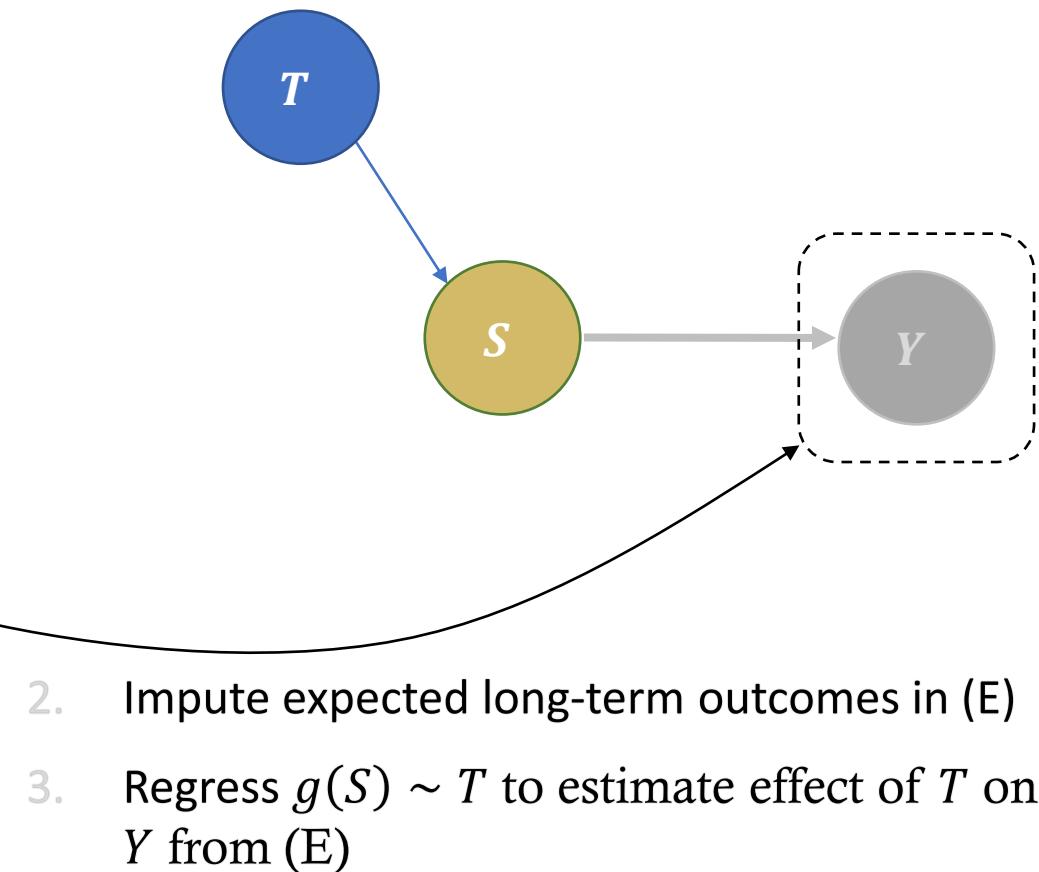
Causal Inference with Surrogates 101

historical/long-term (O)



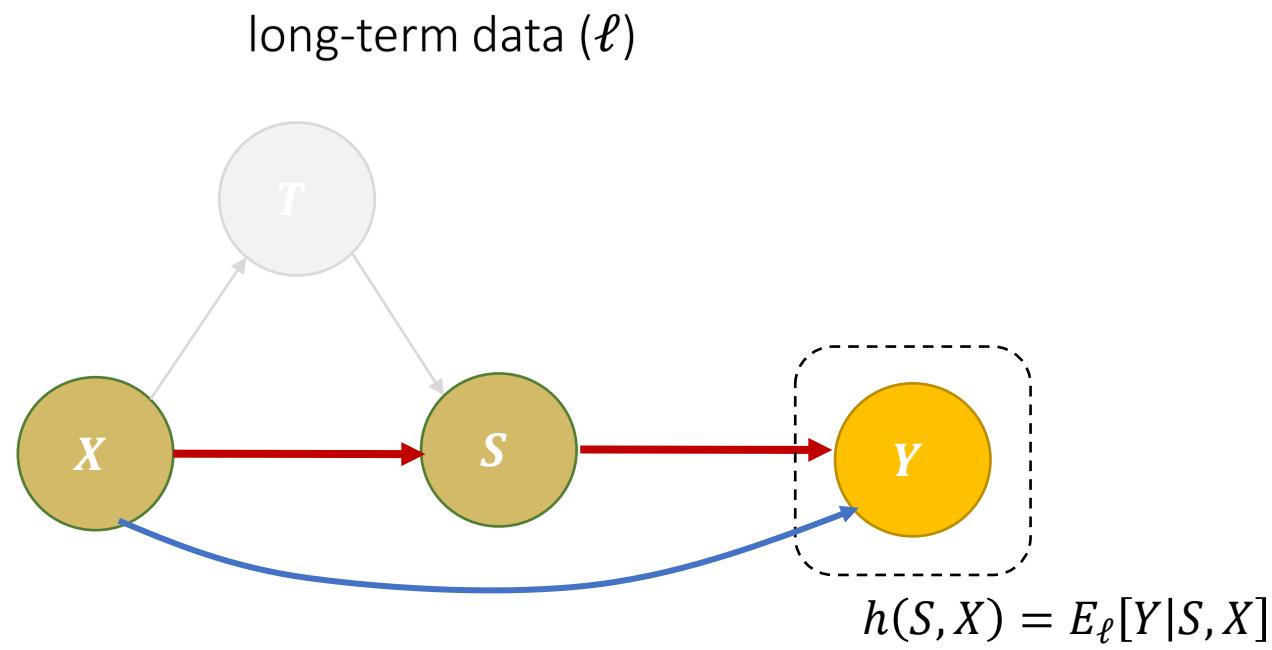
1. Estimate $g(S) := E[Y|S]$ (surrogate index) from (O) by regressing $Y \sim S$

recent/short-term (E)

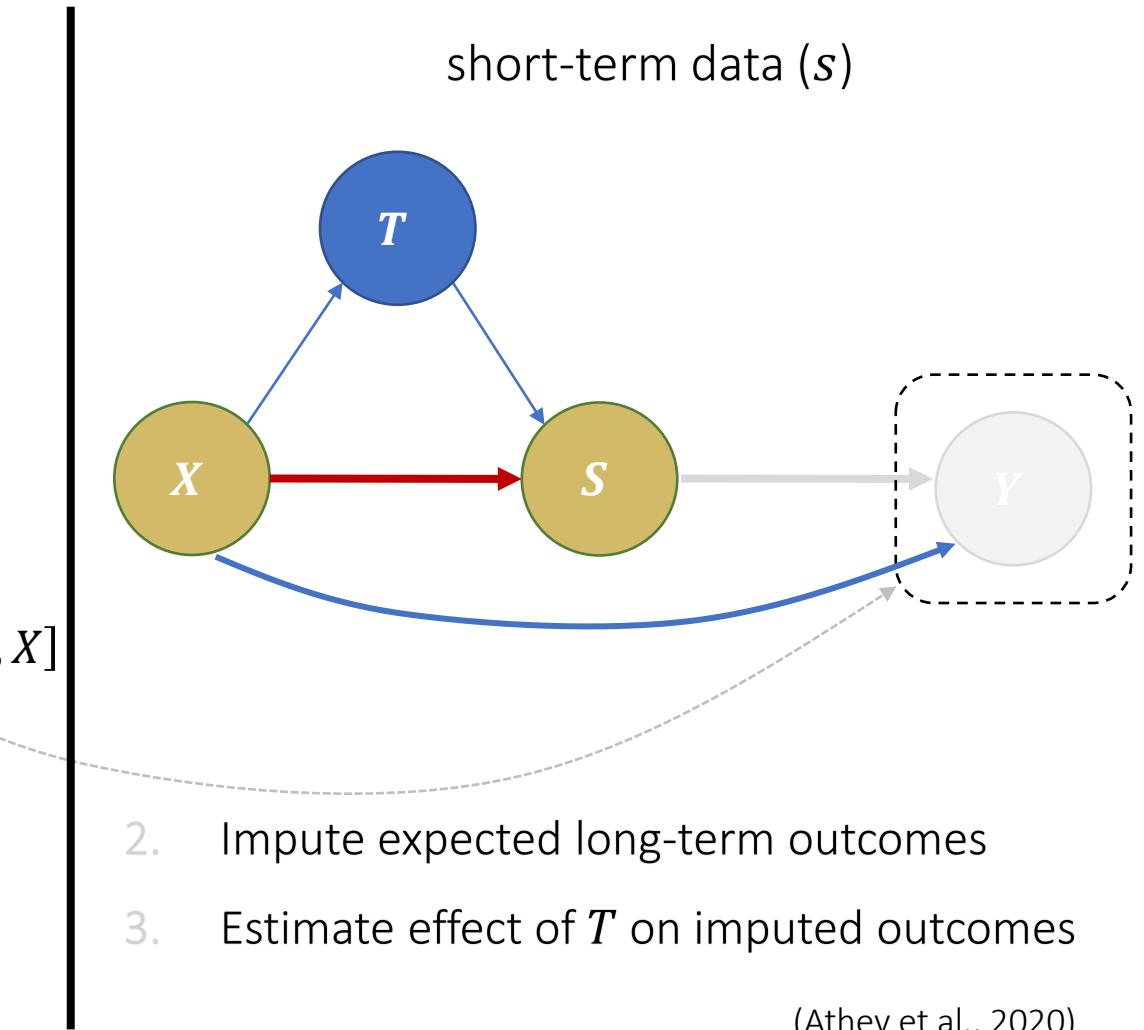


2. Impute expected long-term outcomes in (E)
3. Regress $g(S) \sim T$ to estimate effect of T on Y from (E)

Causal Inference with Surrogates



1. Estimate $h(S, X) := E[Y|S, X]$ (surrogate index)
by regressing $Y \sim S, X$



Causal Inference with Surrogates 102

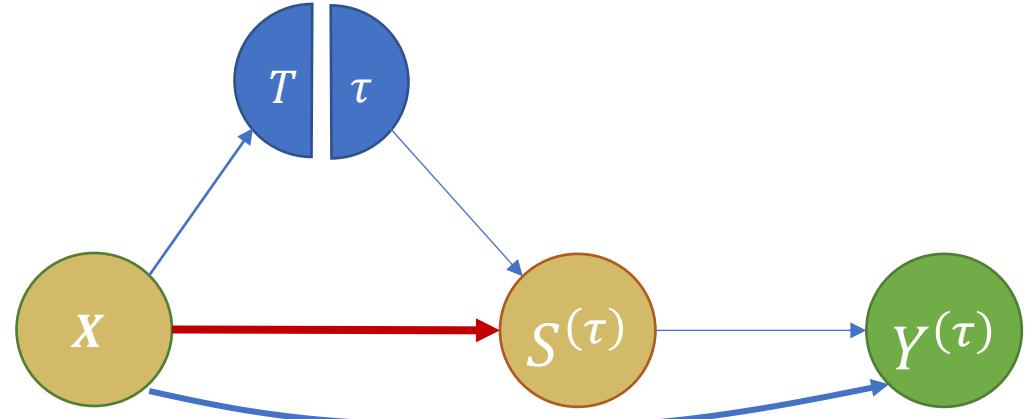
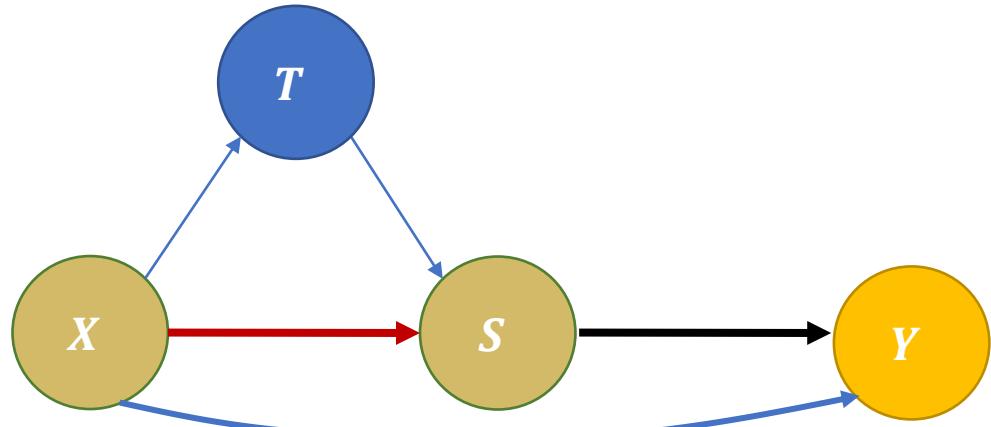
Reminder: $\theta_0 := E[Y^{(\tau)}] = E[E[Y|T = \tau, X]]$, (g-formula)

- Surrogacy assumption implies:

$$E[Y|T = \tau, X] = E[E[Y|S, X] | T = \tau, X]$$

- Surrogacy based g-formula

$$\theta_0 = E[E[h(S, X)|T = \tau, X]]$$



Nested Regression Estimand

- Surrogacy based g-formula

$$\begin{aligned}\theta_0 &= E_s[g(\tau, X)] \\ g(T, X) &= E_s[h(S, X)|T, X] \\ h(S, X) &= E_\ell[Y|S, X]\end{aligned}$$

- g regression function is trained on the outputs of the h regression function
- Surrogate Index Model h is trained on long-term data distribution
- Target parameter is averaged over short-term data distribution

Automatic Debiased Machine Learning

- Surrogacy based g-formula

$$\begin{aligned}\theta_0 &= E_s[g(\tau, X)] \\ g(T, X) &= E_s[h(S, X)|T, X] \\ h(S, X) &= E_\ell[Y|S, X]\end{aligned}$$

- Automatic debiased formula

$$\theta_0 = E_s \left[g(\tau, X) + a_1(T, X) (h(S, X) - g(T, X)) \right]$$

Recursive Automatic Debiasing

- Surrogacy based g-formula

$$\begin{aligned}\theta_0 &= E_s[g(\tau, X)] \\ g(T, X) &= E_s[h(S, X) | T, X] \\ h(S, X) &= E_\ell[Y | S, X]\end{aligned}$$

- Automatic debiased formula

$$\theta_0 = E_s[g(\tau, X) + a_1(T, X)(h(S, X) - g(D, X))] + E_\ell[a_2(S, X)(Y - h(S, X))]$$

Inverse Propensity ratio
 $\frac{1\{T = \tau\}}{\Pr(T = \tau | X)}$

Surrogate Score(S, X) * $E[a_1(T, X) | S, X]$
Surrogate score: (density ratio of S, X between ℓ and s “environment”)

$$\frac{p(S, X | E = s)}{p(S, X | E = \ell)} = \frac{\Pr(E = s | S, X) \Pr(E = \ell)}{\Pr(E = \ell | S, X) \Pr(E = s)}$$

Multiplied by

$$E[a_1(T, X) | S, X] = \frac{\Pr(T = t | S, X)}{\Pr(T = t | X)}$$

An Application from Operations Management

Return on Investment at Microsoft using Surrogates and Dynamic Treatment Effects

Estimating Long-Term Returns on Investment

- Companies frequently deploys new discount or customer support programs
- Which of these programs (“investments”) are more successful than others?
- Success is a **long-term** objective: what is the effect of the program on the two-year customer journey (e.g., effect on two-year revenue)
- We cannot wait two years to evaluate a program
- **Main Question.** Can we construct estimates of the values of these programs with **short-term** data, e.g. after 6 months?

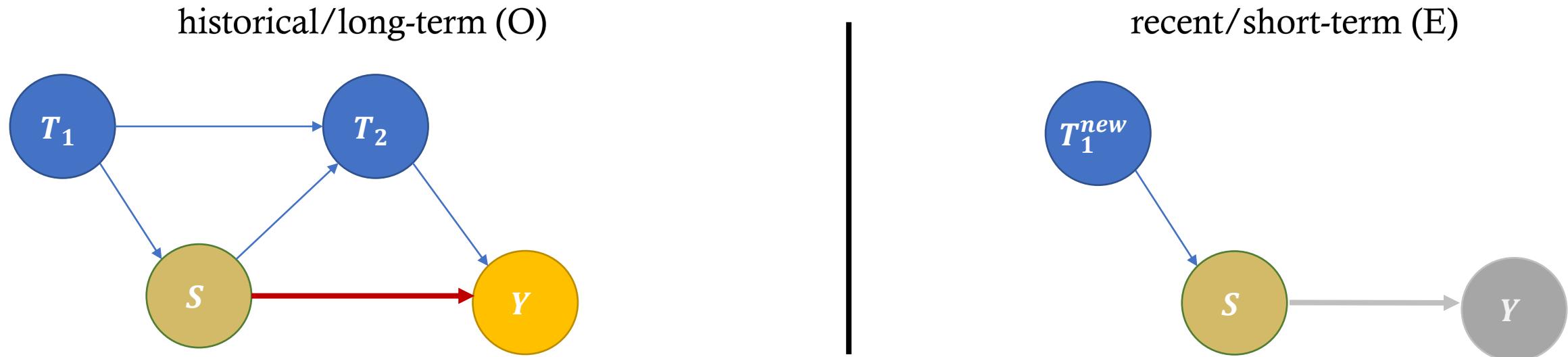


Key Assumptions

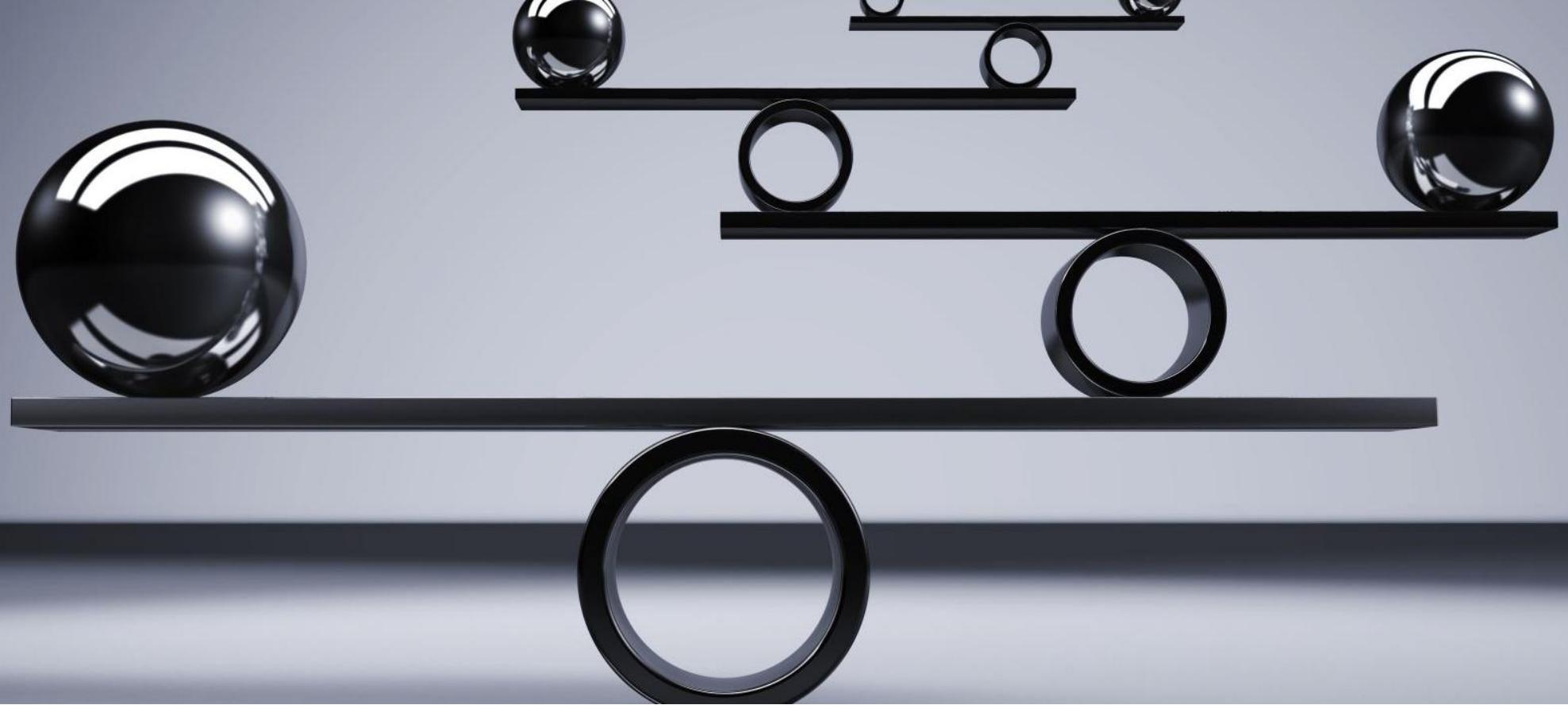
- Long-term effect only goes through surrogates
- Expected relationship between surrogates and long-term reward is the same long-term setting (O) and in short-term setting (E)

Key Assumptions can be Easily Violated

Investment policies are dynamic and change



- ◊ We deployed **older/deprecated** investments
- ◊ In a potentially long-term highly **auto-correlated** manner
- ◊ Investments are potentially **adaptive**
- ◊ Investment policies **change**

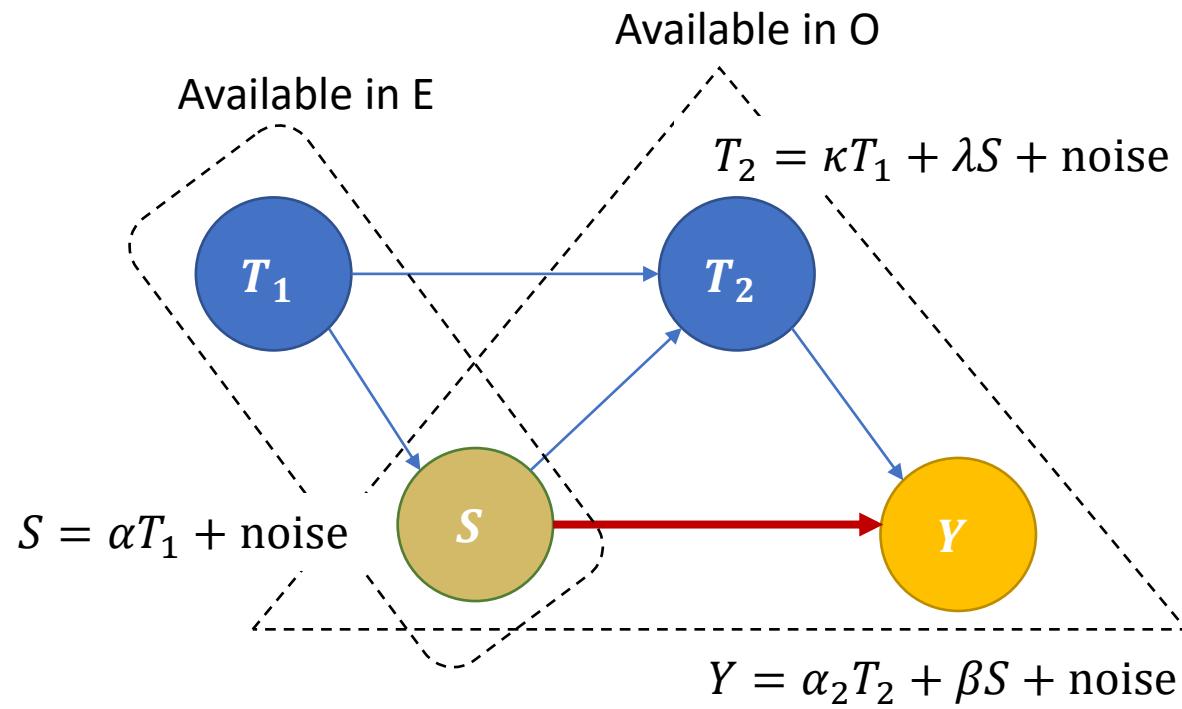


Bias of Vanilla Surrogate

Illustrative Example

Illustrative Example

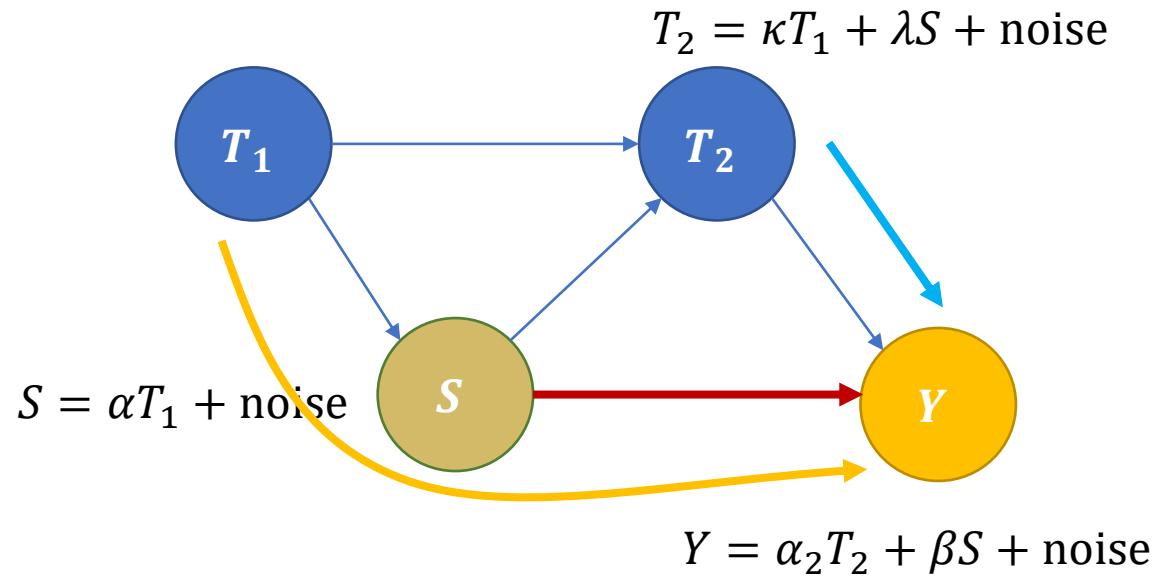
- Want the effect of T_1 on Y , with no future treatments



Bias of Vanilla Surrogate

- ❖ What do we want to estimate?

$$Y = \beta \underset{\substack{!! \\ \theta_0}}{\alpha} T_1 + \alpha_2 T_2 + \text{noise}$$



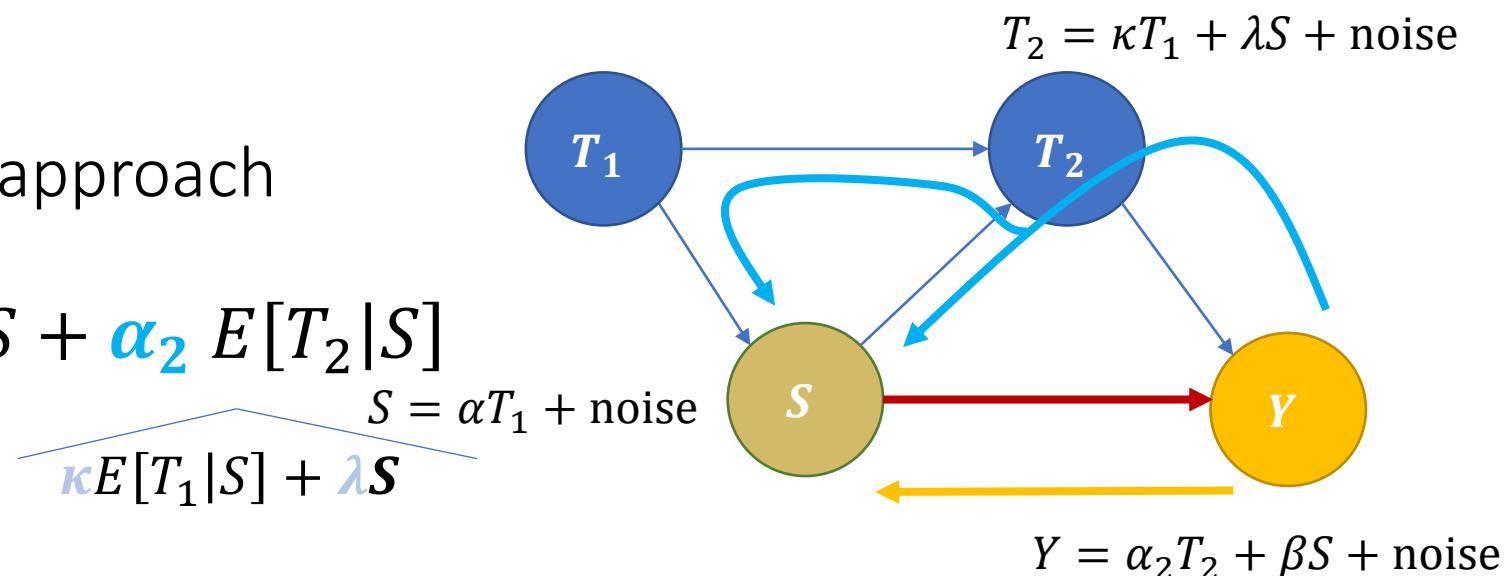
Bias of Vanilla Surrogate

- What do we want to estimate?

$$Y = \beta \alpha_{\theta_0} T_1 + \alpha_2 T_2 + \text{noise}$$

- What does the surrogate approach estimate?

$$g_0(S) := E[Y|S] = \beta S + \alpha_2 E[T_2|S]$$
$$S = \alpha T_1 + \text{noise}$$
$$\kappa E[T_1|S] + \lambda S$$



Bias of Vanilla Surrogate

- What do we want to estimate?

$$Y = \beta \alpha T_1 + \alpha_2 T_2 + \text{noise}$$

!!
 θ_0

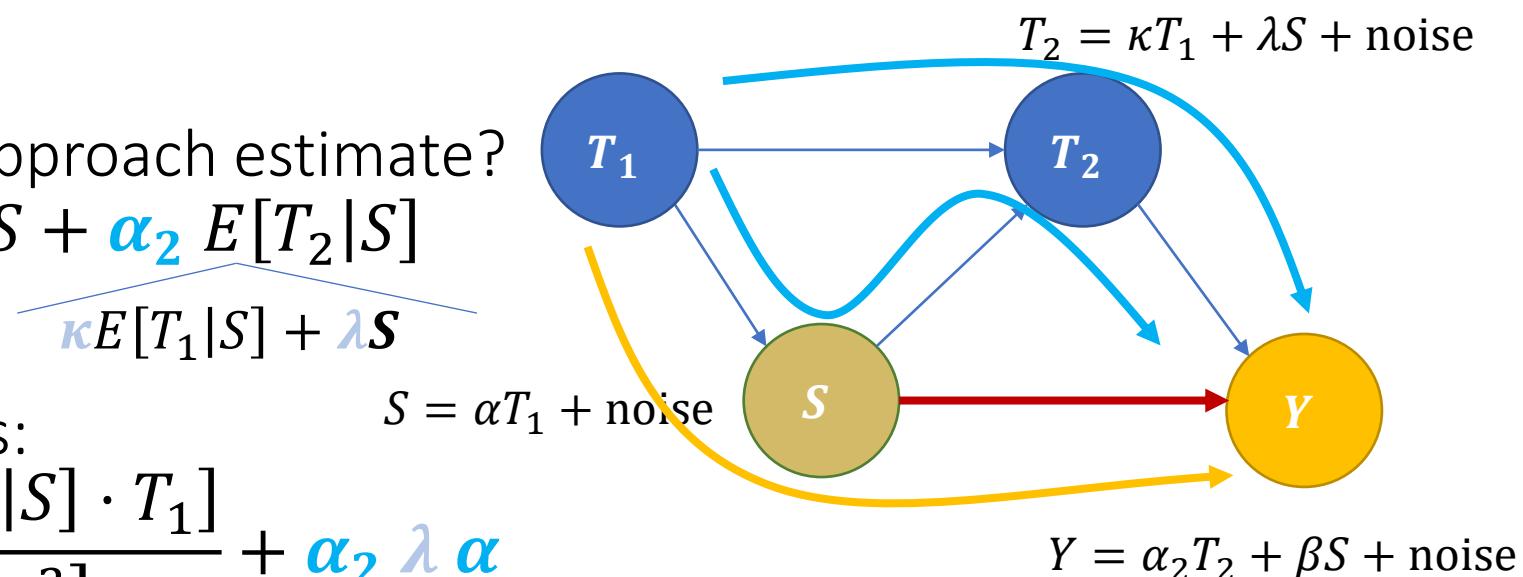
- What does the surrogate approach estimate?

$$g_0(S) := E[Y|S] = \beta S + \alpha_2 E[T_2|S]$$

$\kappa E[T_1|S] + \lambda S$

- The effect of T_1 on $g_0(S)$ is:

$$\theta_* := \beta \alpha + \alpha_2 \kappa \frac{E[E[T_1|S] \cdot T_1]}{E[T_1^2]} + \alpha_2 \lambda \alpha$$



Bias of Vanilla Surrogate

- What do we want to estimate?

$$Y = \beta \underset{\theta_0}{\alpha} T_1 + \alpha_2 T_2 + \text{noise}$$

- What does the surrogate approach estimate?

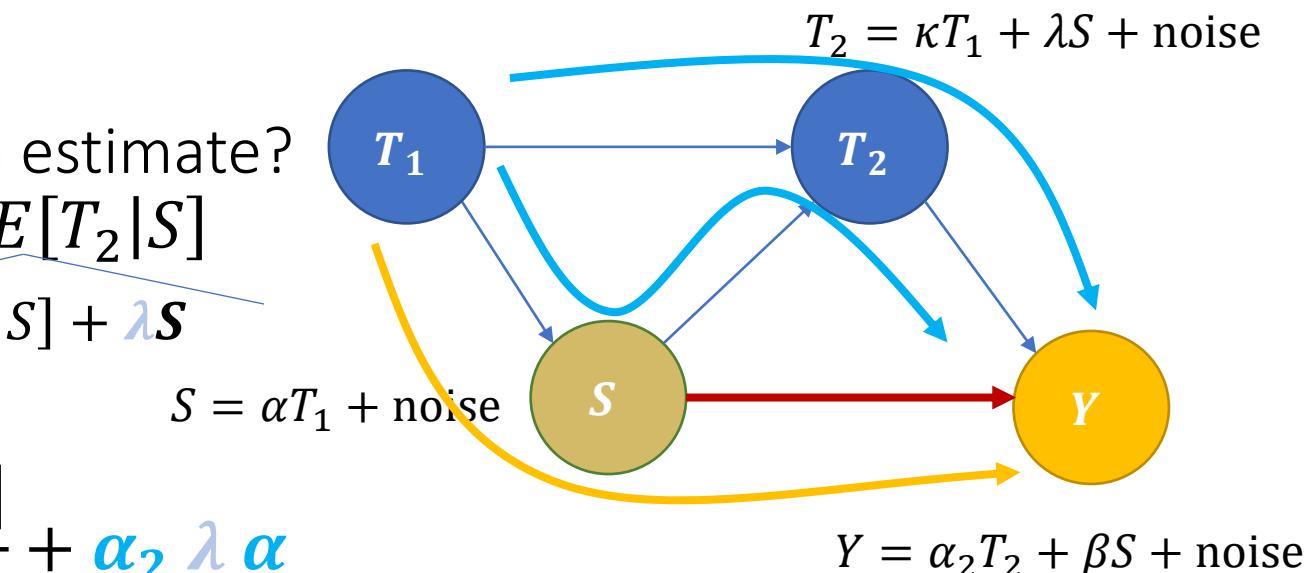
$$g_0(S) := E[Y|S] = \beta S + \alpha_2 E[T_2|S]$$

$$\kappa E[T_1|S] + \lambda S$$

- The effect of T_1 on $g_0(S)$ is:

$$\theta_* := \beta \alpha + \alpha_2 \kappa \frac{E[E[T_1|S] \cdot T_1]}{E[T_1^2]} + \alpha_2 \lambda \alpha$$

Bias from the fact that policy is auto-correlated violating the “effect only through surrogates” assumption



Bias from the fact that policy is adaptive on past surrogates violating “relationship of S and Y unchanged” assumption

Prevalence of Such Violation

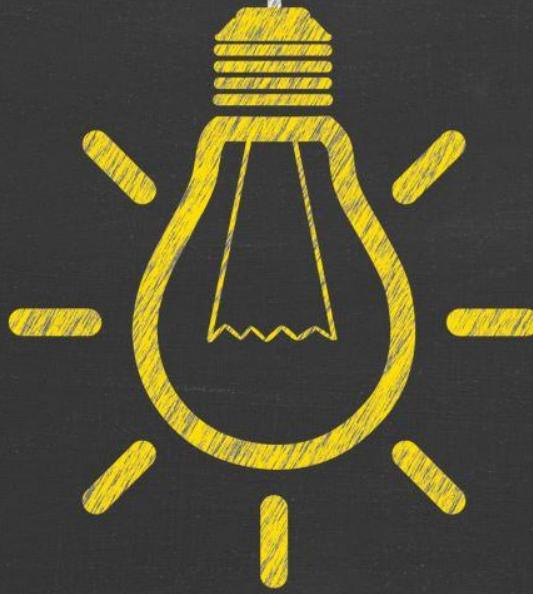
- Marketing departments frequently deploy new marketing campaigns
 - In historical data, older campaigns were deployed
 - Due to targeting, if a customer received an ad, most probably will receive one in the future
 - Ad display is adaptive to signals of customer behavior
- Pharmaceutical companies frequently deploy new drugs
 - Can we get estimates of long-term effect from short-term trials
 - Use short-term signals of patient response that correlated with long-term survival
 - In historical long-term clinical data, patients are treated with multiple treatments over time
 - Treatments are typically adaptive to signals of patient trajectory

Key Idea 1

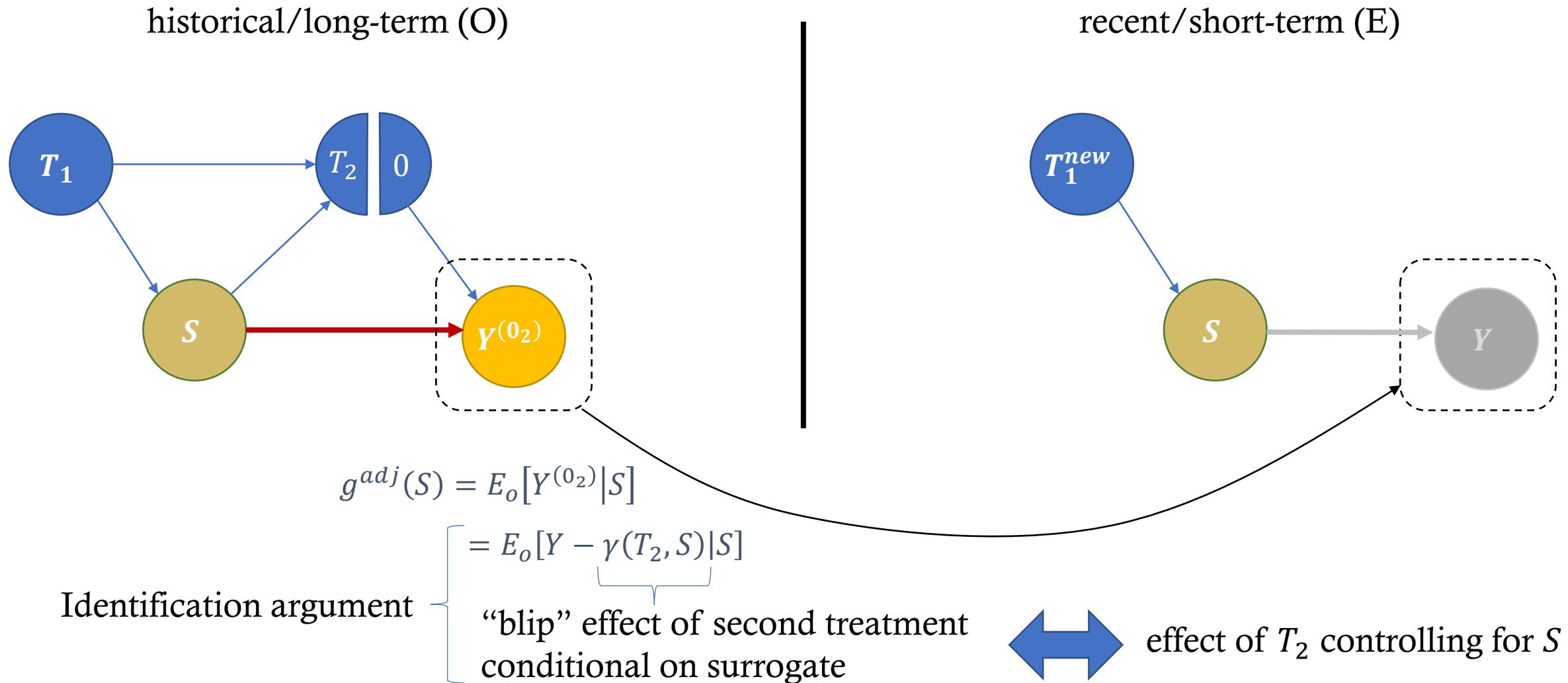
Dynamic Adjustment of Surrogate Index

Battocchi, Dillon, Hei, Lewis, Oprescu, Syrgkanis,
Estimating the Long-Term Effects of Novel
Treatments, NeurIPS'21

<https://arxiv.org/abs/2103.08390>



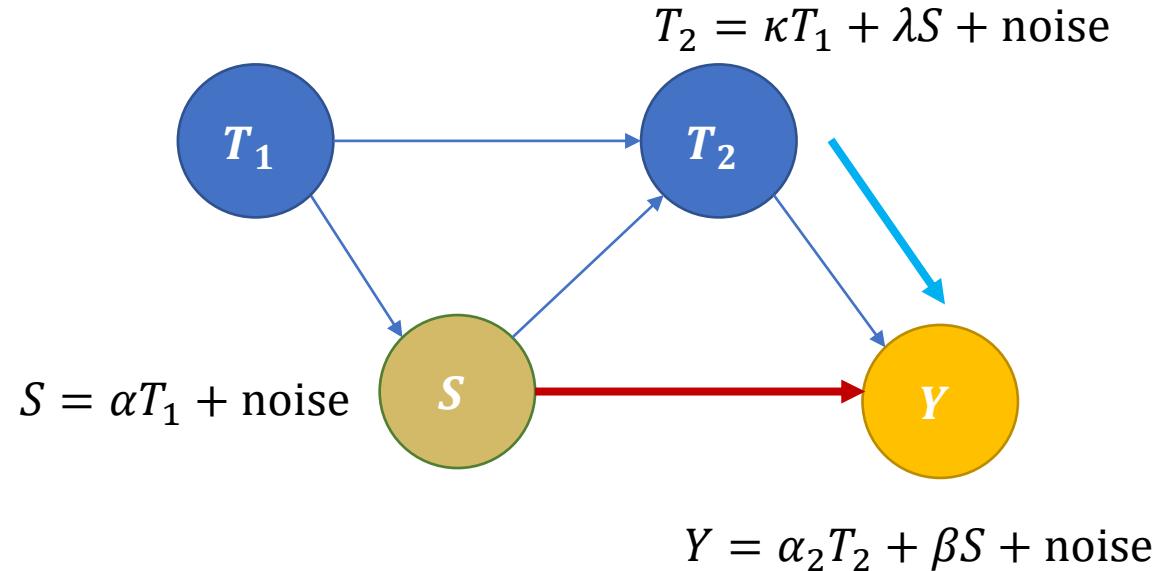
Dynamically Adjusted Surrogate Index



Illustrative Example: Dynamically Adjusted Index

- Estimate effect α_2 of T_2 on Y , controlling for S

$$E[Y|T_2, S] = \alpha_2 T_2 + \beta S$$



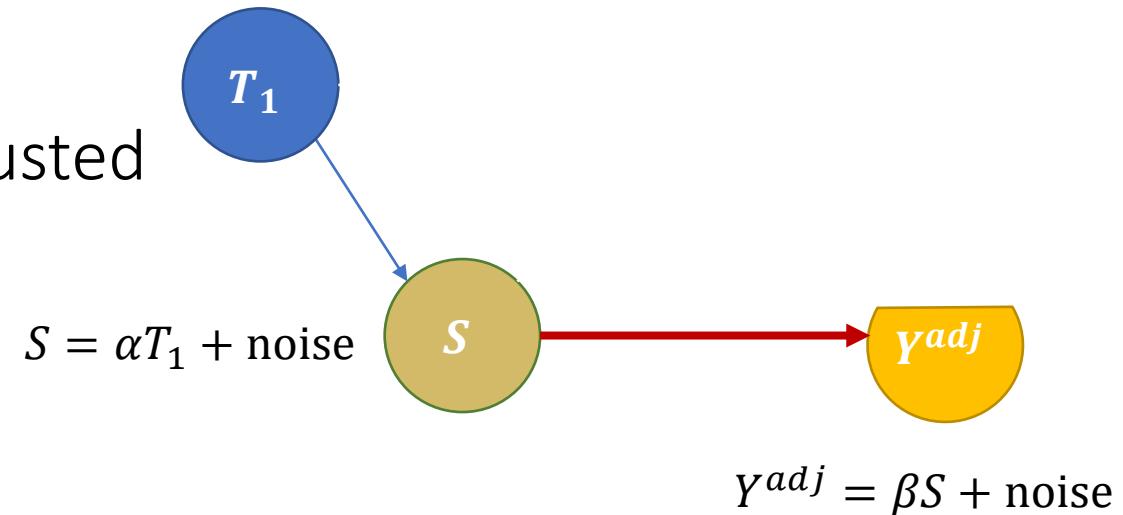
Illustrative Example: Dynamically Adjusted Index

- Estimate effect α_2 of T_2 on Y , controlling for S

$$E[Y|T_2, S] = \alpha_2 T_2 + \beta S$$

- Subtract that effect to create the adjusted outcome

$$Y^{adj} := Y - \alpha_2 T_2$$



Illustrative Example: Dynamically Adjusted Index

- Estimate effect α_2 of T_2 on Y , controlling for S

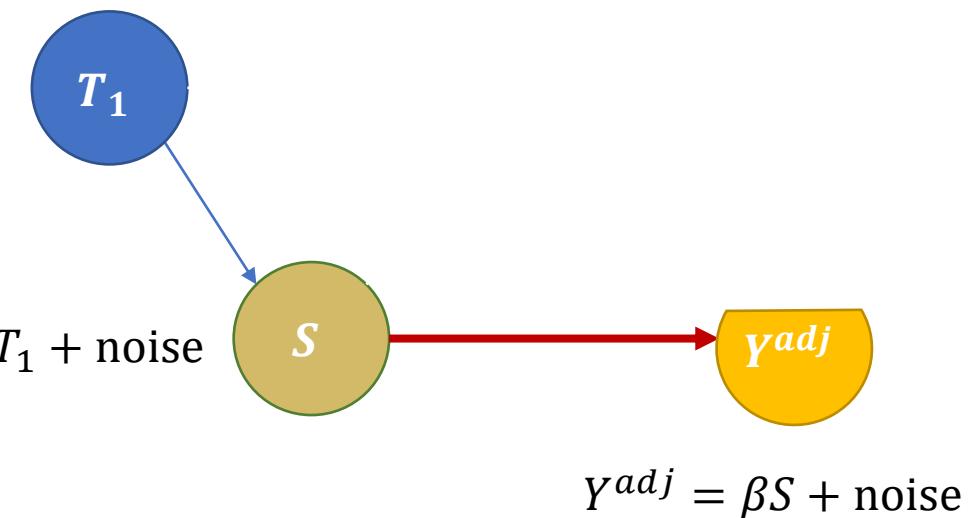
$$E[Y|T_2, S] = \alpha_2 T_2 + \beta S$$

- Subtract that effect to create the adjusted outcome

$$Y^{adj} := Y - \alpha_2 T_2$$

- Estimate dynamically adjusted index on long-term data

$$g_0^{adj}(S) := E[Y^{adj}|S] = \beta S$$



Illustrative Example: Dynamically Adjusted Index

- Estimate effect α_2 of T_2 on Y , controlling for S

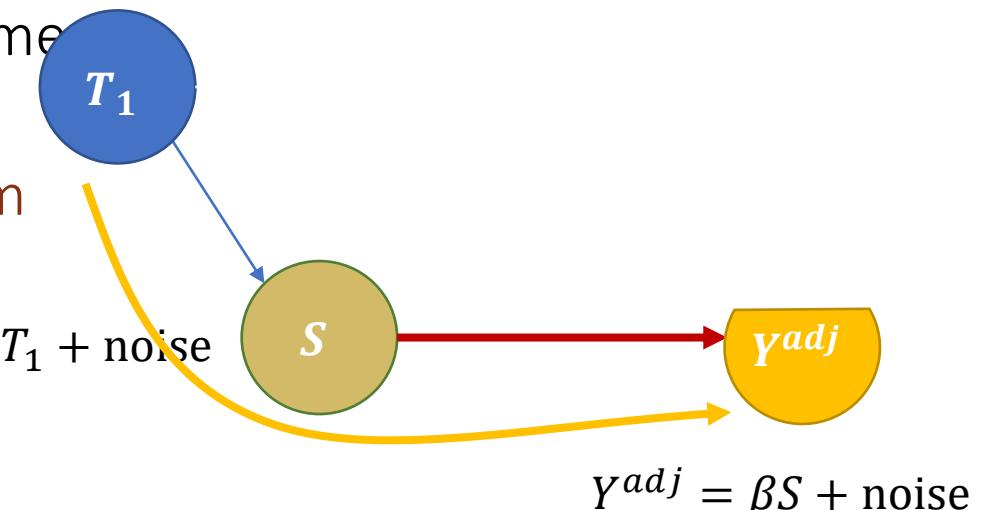
$$E[Y|T_2, S] = \alpha_2 T_2 + \beta S$$

- Subtract that effect to create the adjusted outcome

$$Y^{adj} := Y - \alpha_2 T_2$$

- Estimate dynamically adjusted index **on long-term data**

$$g_0^{adj}(S) := E[Y^{adj}|S] = \beta S$$

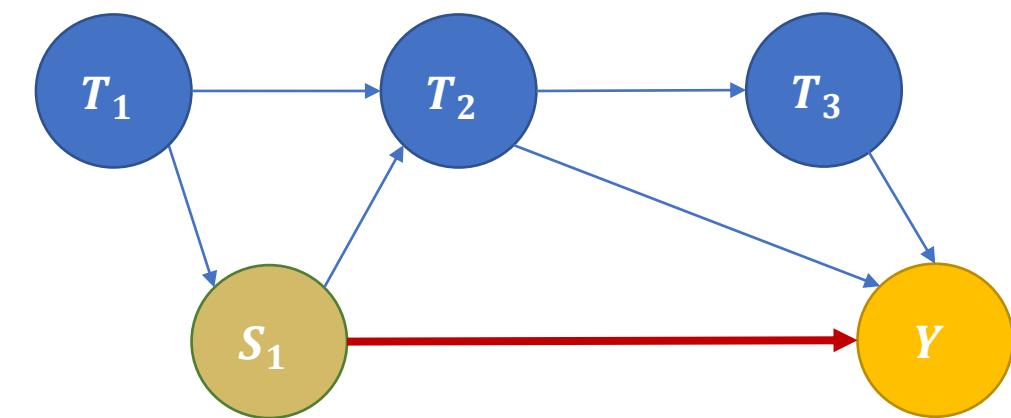


- Estimate effect of T_1 on $g_0^{adj}(S)$ **on short-term data**

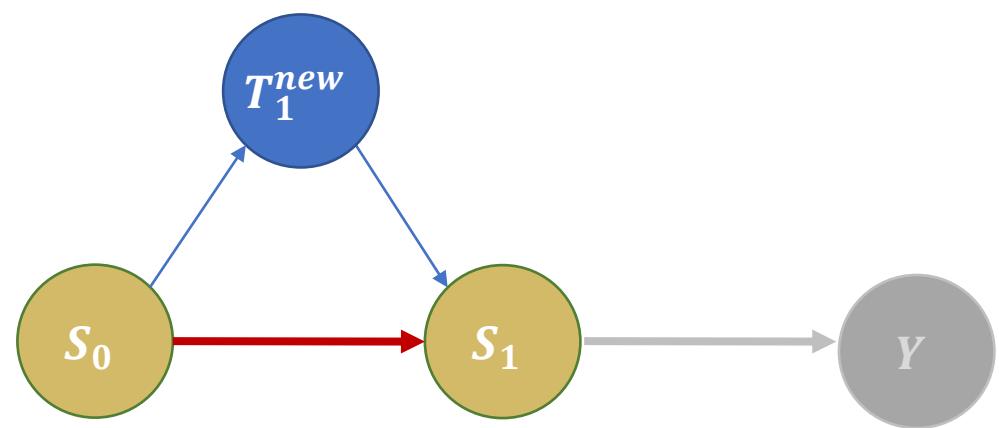
$$\theta_0 := \frac{E[g_0^{adj}(S) \cdot T_1]}{E[T_1^2]} = \beta \cdot \alpha$$

Beyond Two Periods

historical/long-term (O)



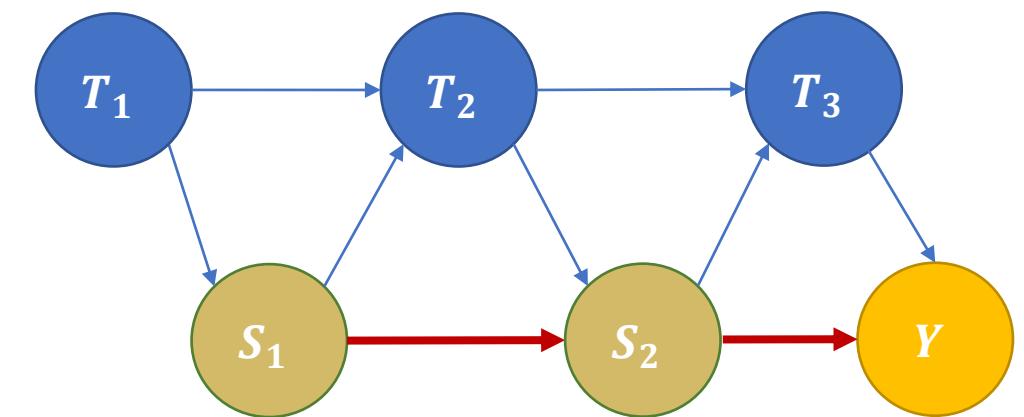
recent/short-term (E)



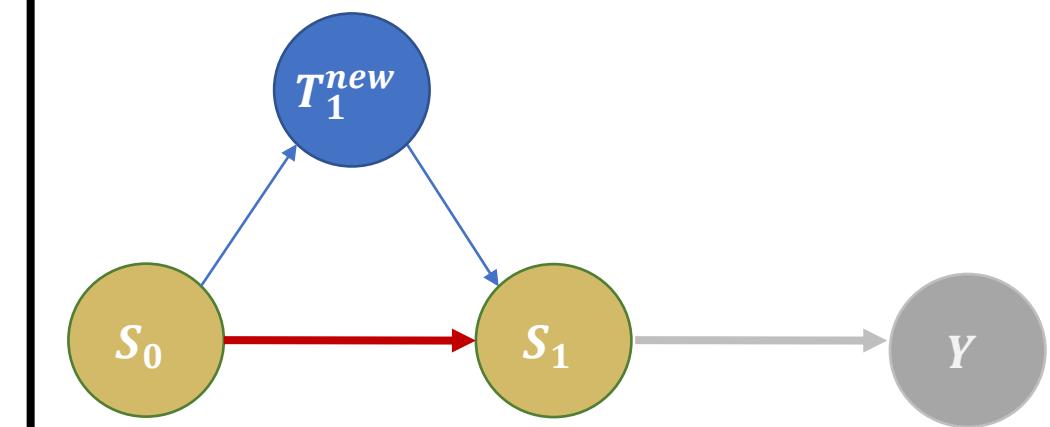
- ❖ Multiple treatments are offered after the surrogate variable
- ❖ Easy: estimate their effect controlling for S_1
- ❖ Wrong!

Beyond Two Periods

historical/long-term (O)

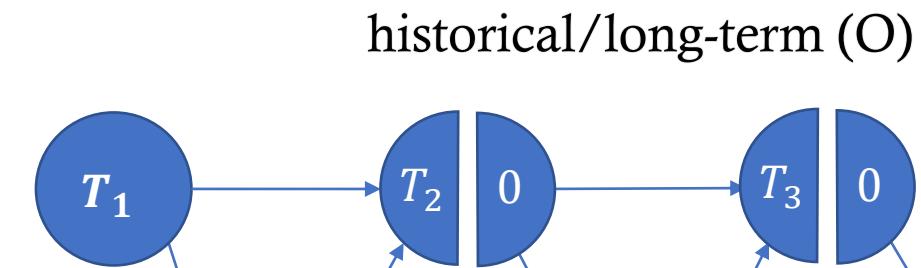


recent/short-term (E)



- ❖ Treatments are offered in an adaptive manner, in response to previous period surrogate/state
- ❖ The surrogate – treatment feedback precludes viewing this as a one-shot treatment problem
- ❖ Setting is known as the **dynamic treatment regime** [Robins'94, '04, Chakraborty-Murphy'14]

Beyond Two Periods: Target Quantity

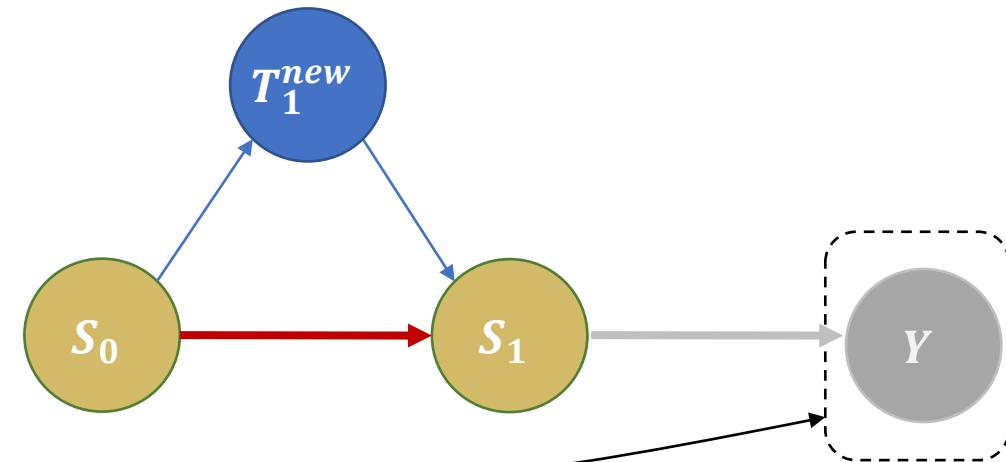


$$g^{adj}(S_1) = E_o[Y^{(0_{\geq 2})} | S_1]$$

$$= E_o \left[Y - \sum_{t \geq 2} \underbrace{\gamma_t(T_t, S_{t-1})}_{\text{“blip” effect of treatment at period } t} \middle| S_1 \right]$$

“blip” effect of treatment at period t

recent/short-term (E)



- ◊ What effect do we subtract from Y ?
- ◊ Do we estimate effect of T_2, T_3 controlling for S_1, S_2 ?
Wrong
- ◊ Do we estimate effect of T_2, T_3 controlling for S_1 ?
Wrong