

# Causal Fairness Analysis

(Causal Inference II - **Lecture 2**)

Elias Bareinboim



Drago Plecko



Columbia University  
Computer Science



# Reference:

D. Plecko, E. Bareinboim.

Causal Fairness Analysis.

TR R-90, CausalAI Lab, Columbia University.

<https://causalai.net/r90.pdf>

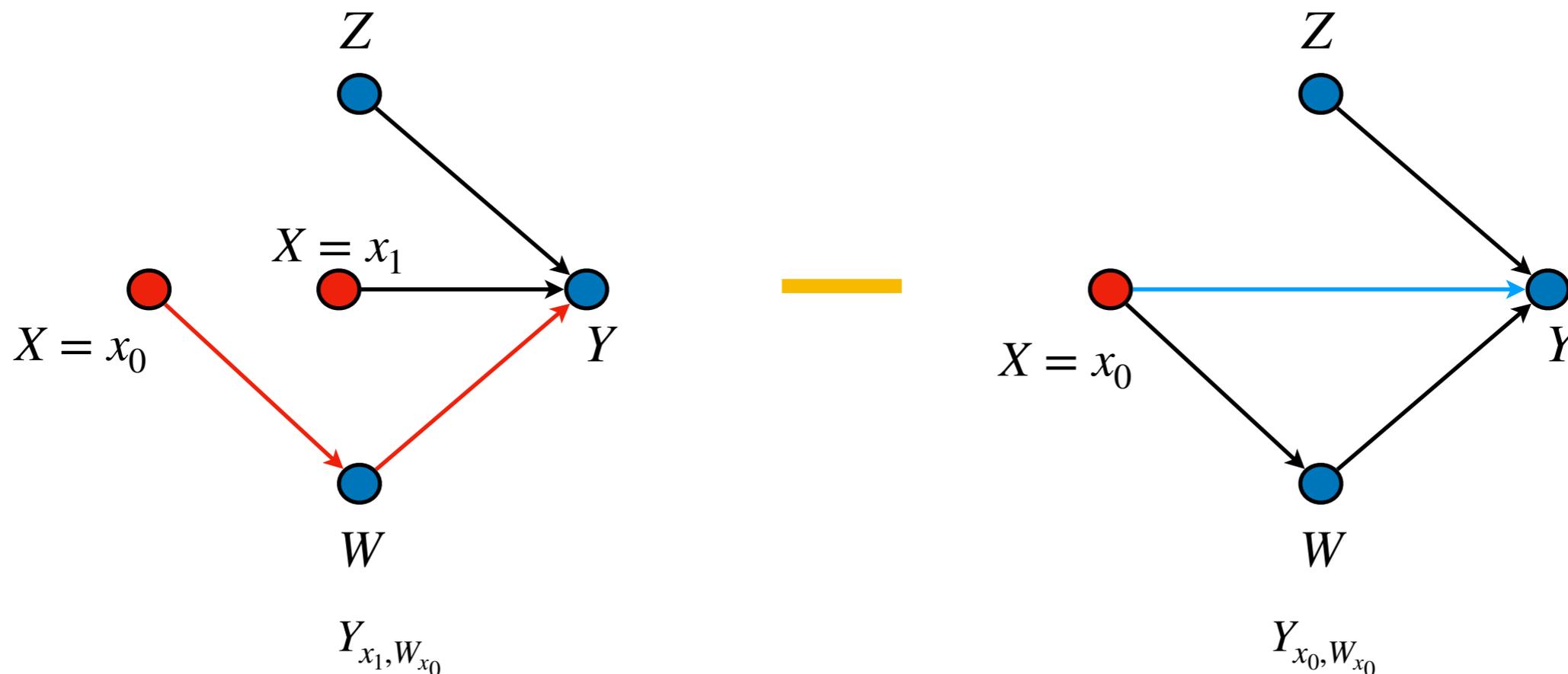
# TV family of causal fairness measures

Section 4

# Gedankenexperiment (NDE)

- For an individual assigned to male ( $X = x_0$ ) by intervention, how would his salary ( $Y$ ) change **had he been** assigned female ( $X = x_1$ ), while keeping the age, nationality, education and employment status unchanged (at the natural level  $X = x_0$ )?

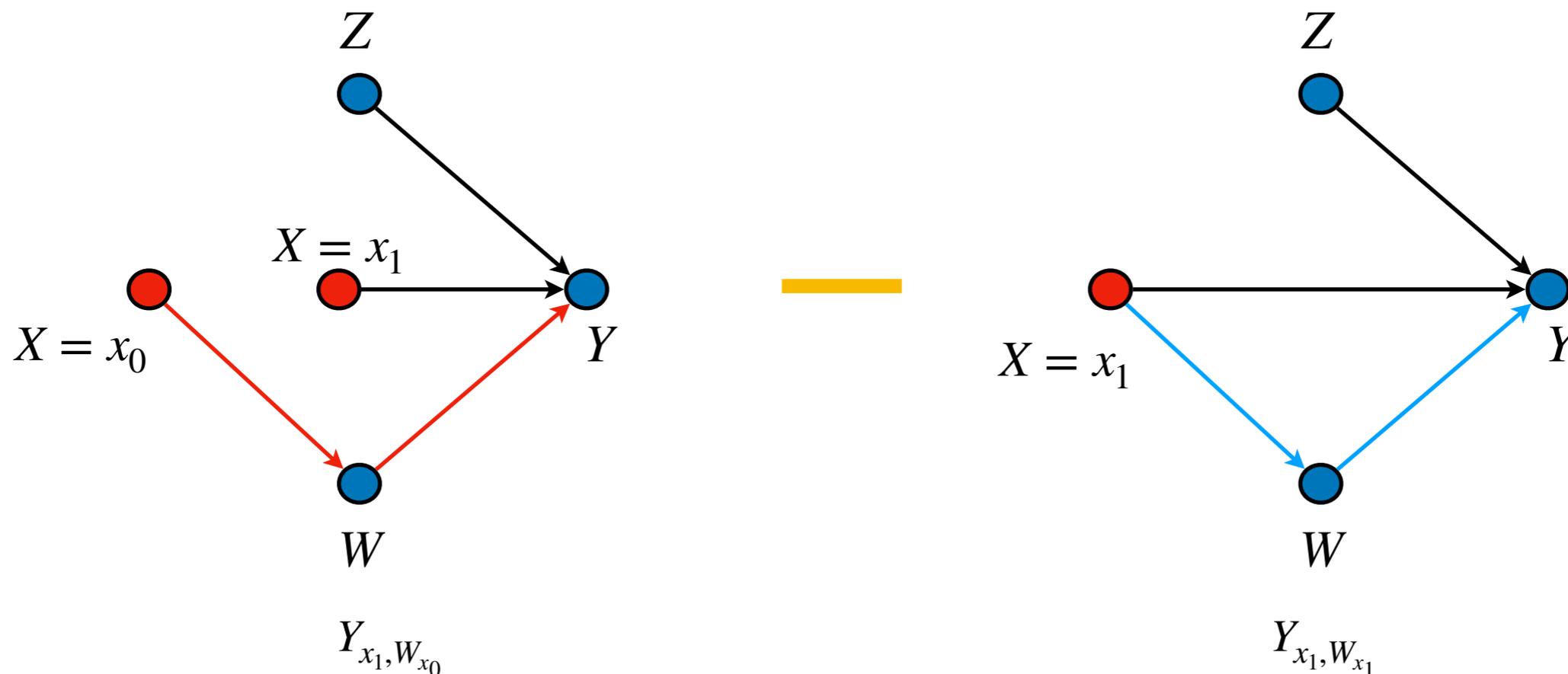
$$\mathbf{NDE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0, W_{x_0}})$$



# Gedankenexperiment (NIE)

- For an individual assigned to be female ( $X = x_1$ ) by intervention, how would her salary ( $Y$ ) change **had she been** assigned to be male ( $X = x_0$ ), while keeping gender unchanged along the direct causal pathway (at the natural level  $X = x_1$ )?

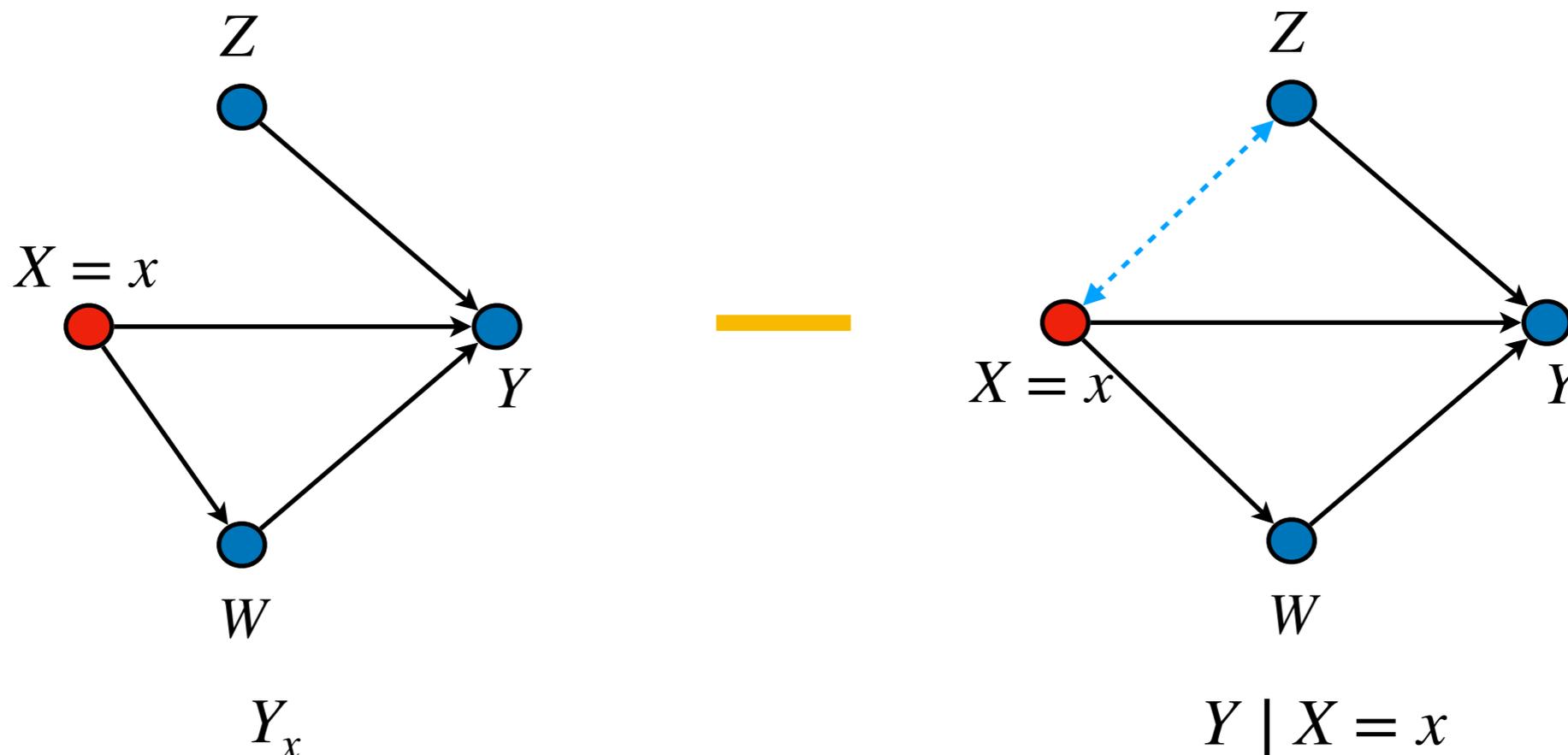
$$\mathbf{NIE}_{x_1, x_0}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_1, W_{x_1}})$$

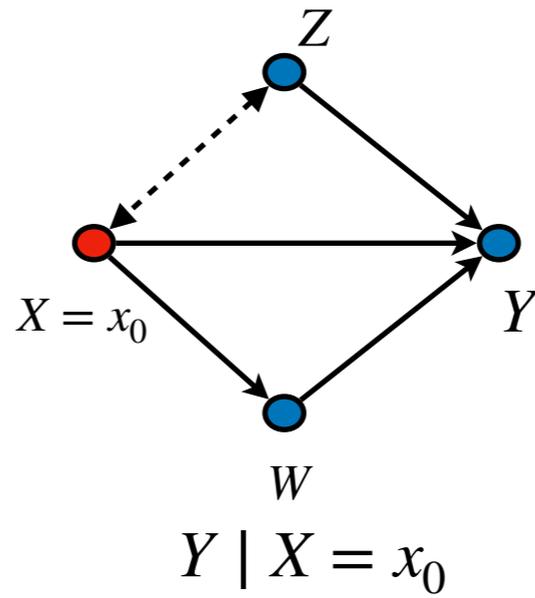
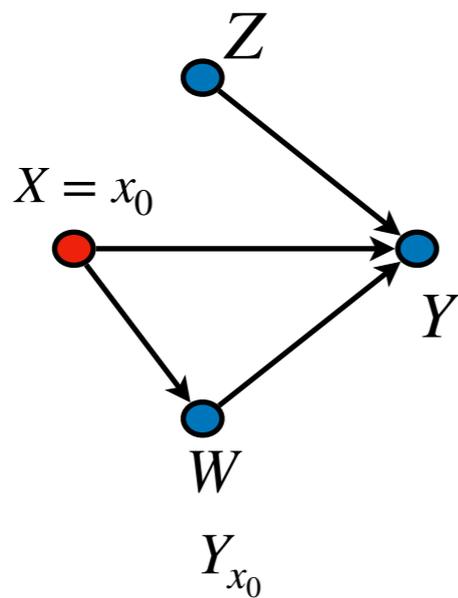
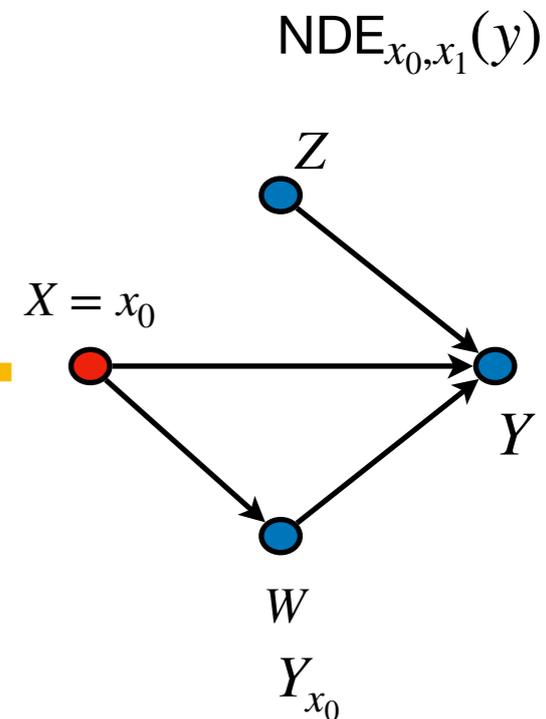
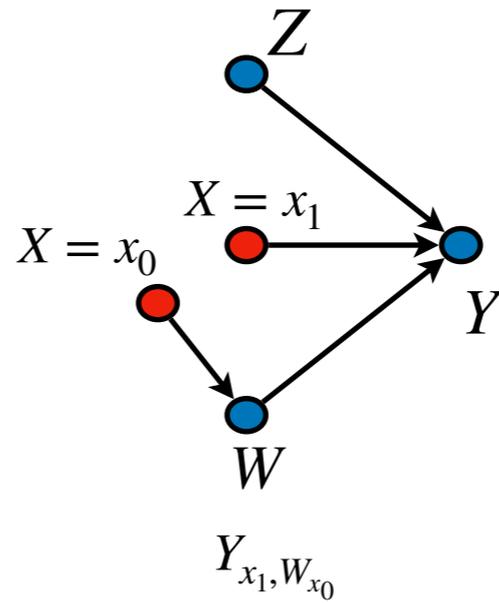
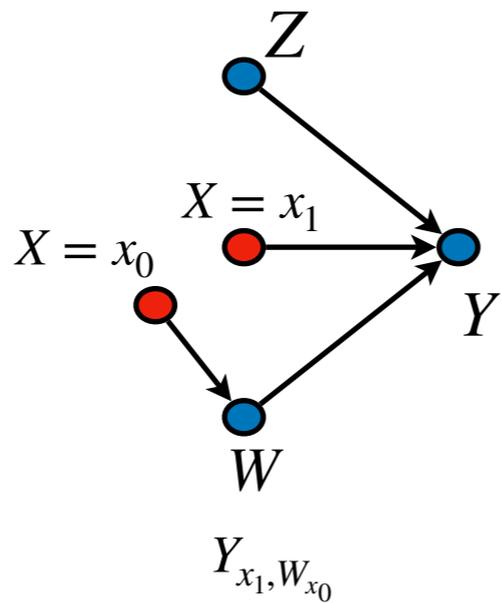
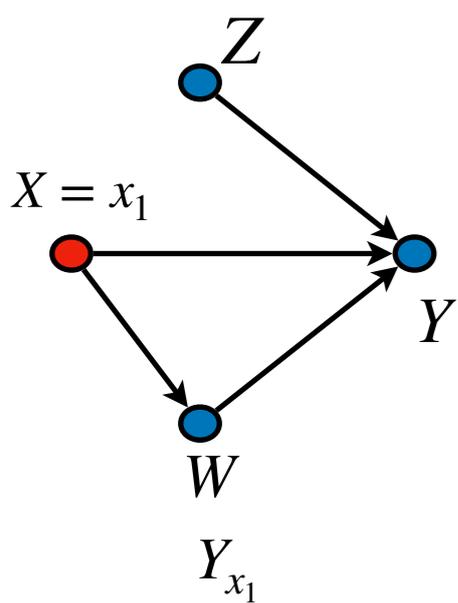
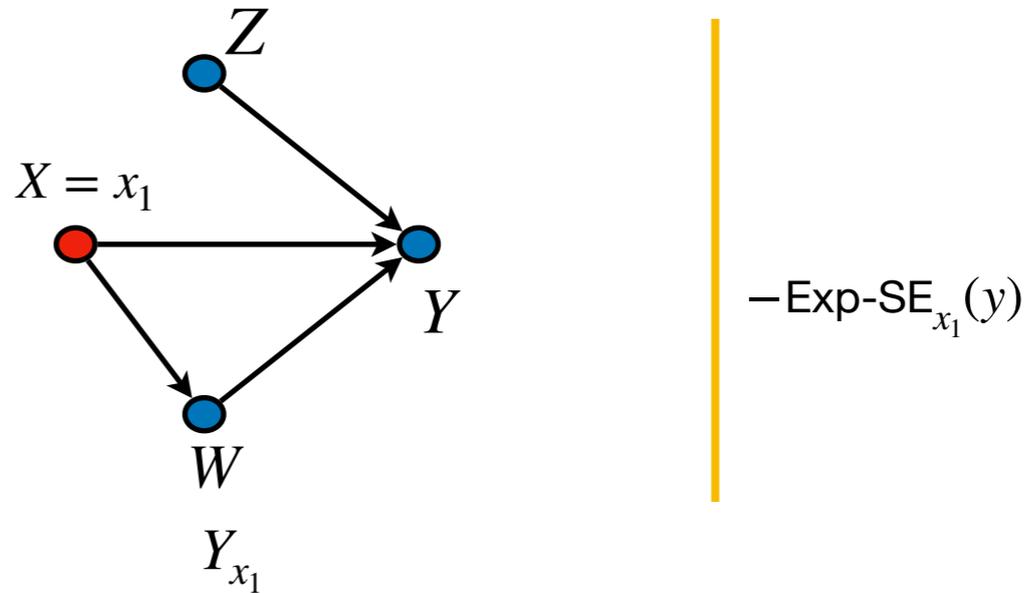
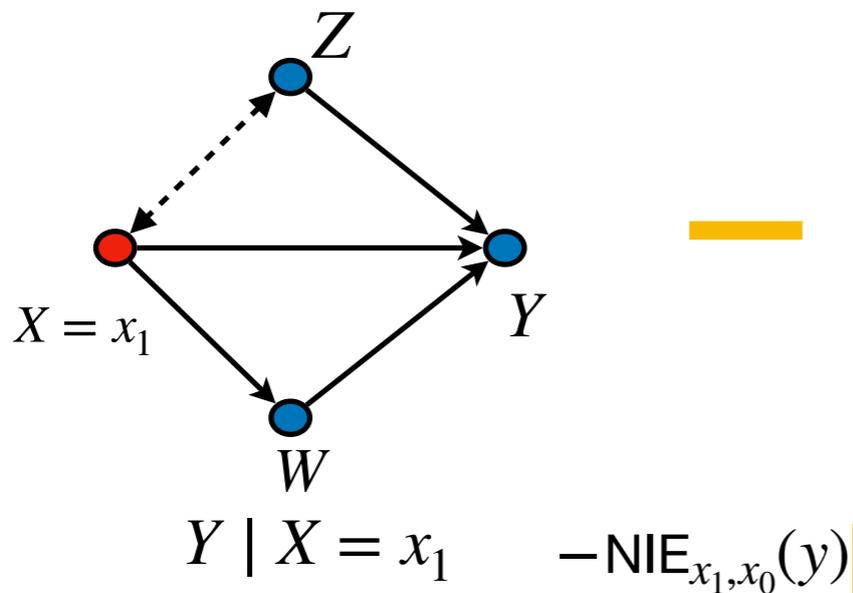


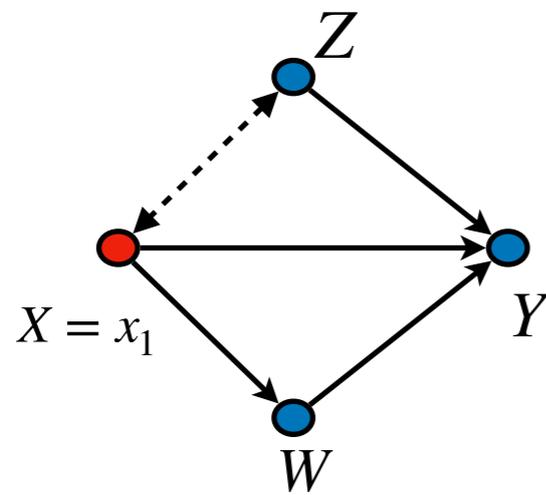
# Gedankenexperiment (Exp-SE)

- How would an individual's salary ( $Y$ ) change if their gender is set to male (or female) by intervention, compared to observing their salary as male (female)?

$$\text{Exp-SE}_x(y) = P(y_x) - P(y | x)$$

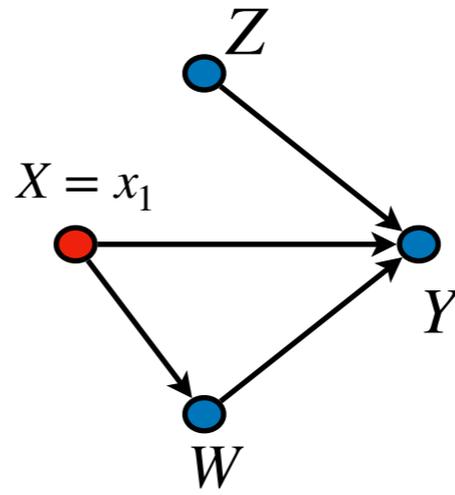






$Y | X = x_1$

$-NIE_{x_1, x_0}(y)$



$Y_{x_1}$

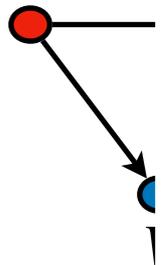
$-\text{Exp-SE}_{x_1}(y)$



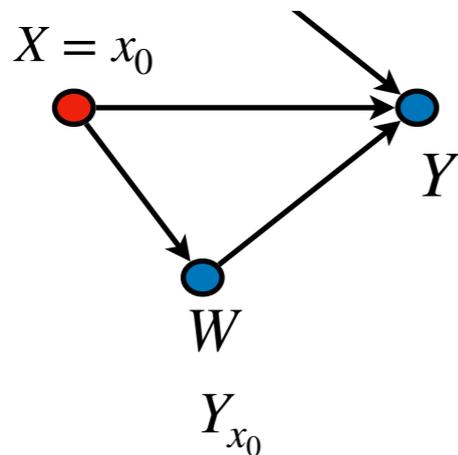
$NDE_{x_0, x_1}(y)$



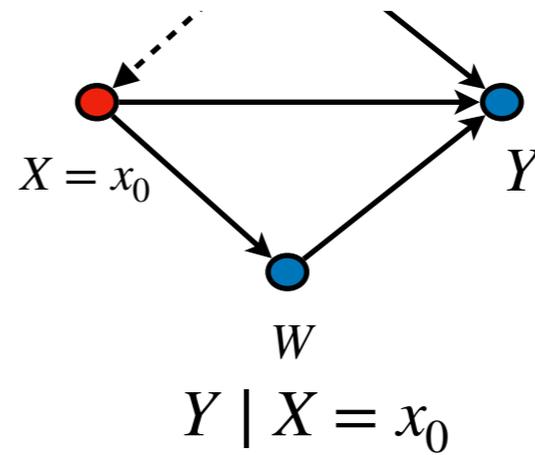
$X = x_1$



# TV Decomposition I



$Y_{x_0}$



$Y | X = x_0$

$\text{Exp-SE}_{x_0}(y)$

# Relation to Structural Fairness

**Corollary.** *The criteria based on NDE, NIE, and Exp-SE measures are **admissible** with respect to structural direct, indirect, and spurious fairness. Formally, these facts are written as:*

$S-DE \implies NDE\text{-fair}$

$S-IE \implies NIE\text{-fair}$

$S-SE \implies Exp\text{-SE-fair}$

admissibility w.r.t.  
structural

In practice, for example, by computing the NDE, we can test for the presence of structural direct effect.

# Testing Structural Fairness in Practice

- Our previous corollary shows that

$$\text{S-DE} \implies \text{NDE-fair} .$$

- By taking this statement's contrapositive, we can see that

$$\text{NDE}_{x_0, x_1}(y) \neq 0 \implies \neg \text{S-DE} .$$

- Therefore, in practice, one may use the following hypothesis testing procedure for testing structural direct effect,

$$H_0 : \text{NDE}_{x_0, x_1}(y) = 0 .$$

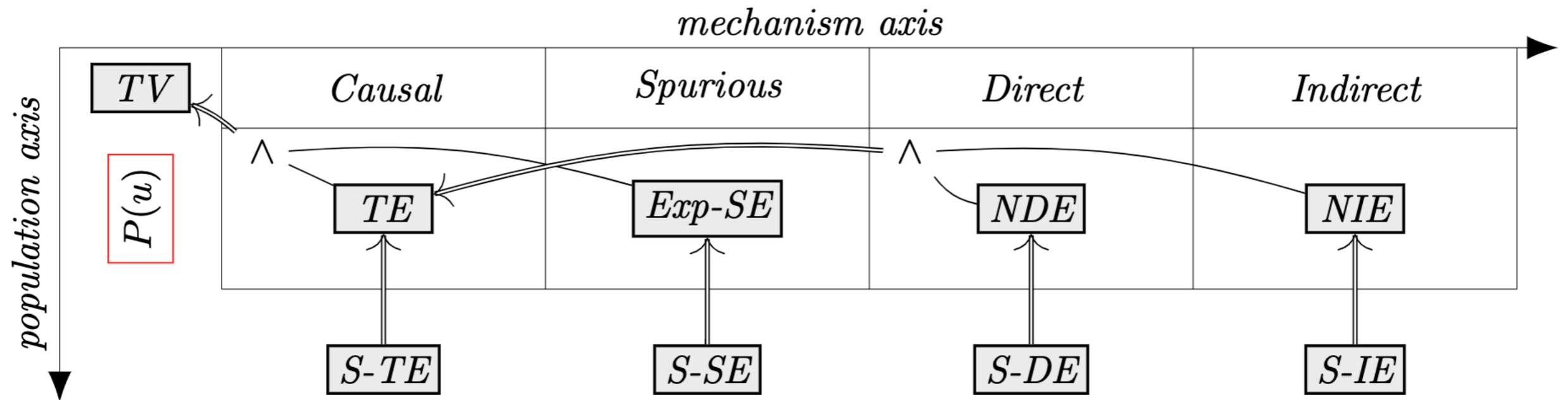
A similar approach can be used for the NIE and Exp-SE since

$$\text{S-IE} \implies \text{NIE-fair}$$

$$\text{S-SE} \implies \text{Exp-SE-fair}$$

This will be used to connect with the disparate treatment and impact doctrines later on.

# Fairness Map (prelim version)



- The map is constructed based on the Corollary in the previous page
- We have found fairness measures that are (i) computable from the data; (ii) admissible with respect to structural fairness; (iii) satisfy decomposability with respect to TV; ✓

Does that mean we are done with Causal Fairness Analysis?

Section 4.2  
Figure 4.2

**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

**SCM  $M^*$**   
(unobserved)

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$\begin{aligned} \text{NDE}_{x_0, x_1}(y) &= P(y_{x_1, W_{x_0}}) - P(y_{x_0}) \\ &= P(\text{Bernoulli}\left(\frac{1}{5}(1 - Z) + \frac{1}{6}W\right) = 1) \\ &\quad - P(\text{Bernoulli}\left(\frac{1}{5}Z + \frac{1}{6}W\right) = 1) \\ &= \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} P(w) \left[ \frac{1}{5}(1 - 2z) + \frac{1}{6}w - \frac{1}{6}w \right] \\ &= \sum_{z \in \{0,1\}} \frac{1}{5}(1 - 2z) = 0. \end{aligned}$$

**Section 4.2**  
**Example 4.1**

**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

**NDE is admissible w.r.t. S-DE.**

**However, here  $NDE = 0$ , and structural direct effect exists.**

**Q: Is NDE powerful enough for detecting discrimination?**

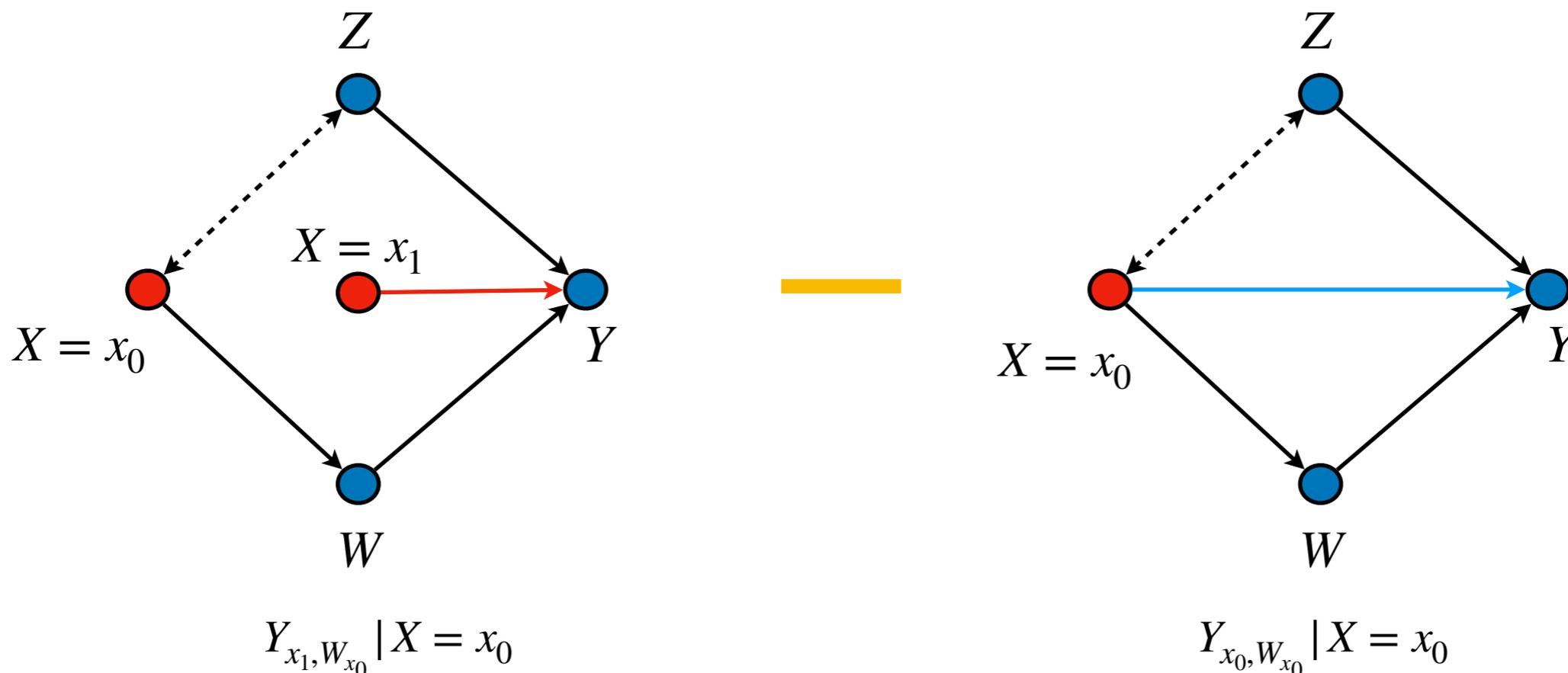
offer  
w]  
4.2

e 4.1

# Gedankenexperiment (Ctf-DE)

- For a male person  $X = x_0$ , how would his salary change ( $Y$ ) **had he been** a female ( $X = x_1$ ), while keeping the age, nationality, education and employment status unchanged (at the level of  $X = x_0$ )?

$$\text{Ctf-DE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}} | x_0) - P(y_{x_0, W_{x_0}} | x_0)$$



**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

### SCM M

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$\begin{aligned} \text{Ctf-DE}_{x_0, x_1}(y | x_0) &= P(y_{x_1, W_{x_0}} | x_0) - P(y_{x_0} | x_0) \\ &= P(\text{Bernoulli}\left(\frac{1}{5}(1 - Z) + \frac{1}{6}W\right) = 1 | x_0) \\ &\quad - P(\text{Bernoulli}\left(\frac{1}{5}Z + \frac{1}{6}W\right) = 1 | x_0) \\ &= \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} P(w)P(z | x_0) \left[ \frac{1}{5}(1 - 2z) + \frac{1}{6}w - \frac{1}{6}w \right] \\ &= \sum_{z \in \{0,1\}} \frac{1}{5}(1 - 2z)P(z | x_0) = 0.036. \end{aligned}$$

**Section 4.2**  
**Example 4.2**

**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

Key properties of Ctf-DE:

1. Ctf-DE is admissible.

2. Ctf-DE is more powerful than NDE.

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(U)$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$= \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} P(w)P(z | x_0) \left[ \frac{1}{5}(1 - 2z) + \frac{1}{6}w - \frac{1}{6}w \right]$$

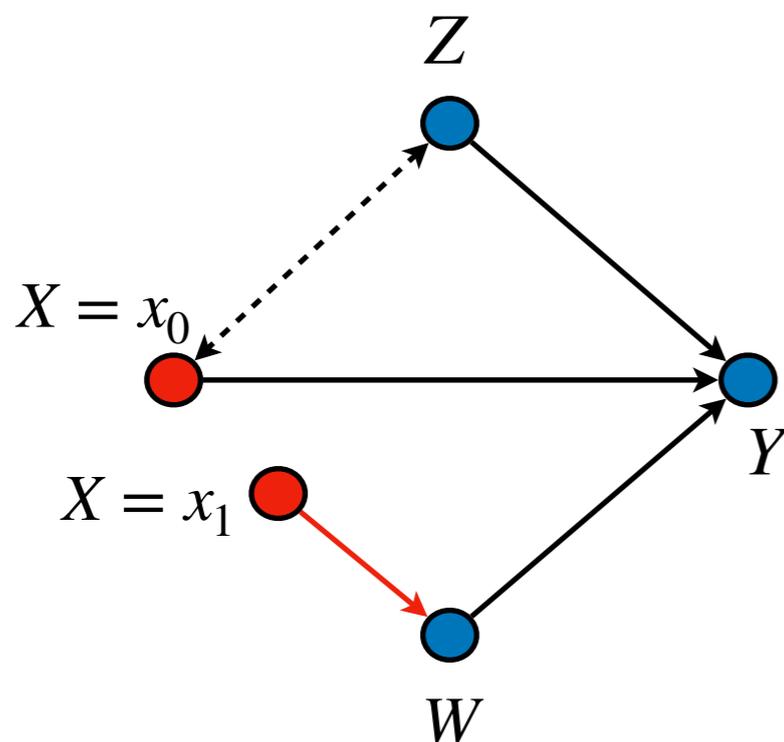
$$= \sum_{z \in \{0,1\}} \frac{1}{5}(1 - 2z)P(z | x_0) = 0.036.$$

Section 4.2  
Example 4.2

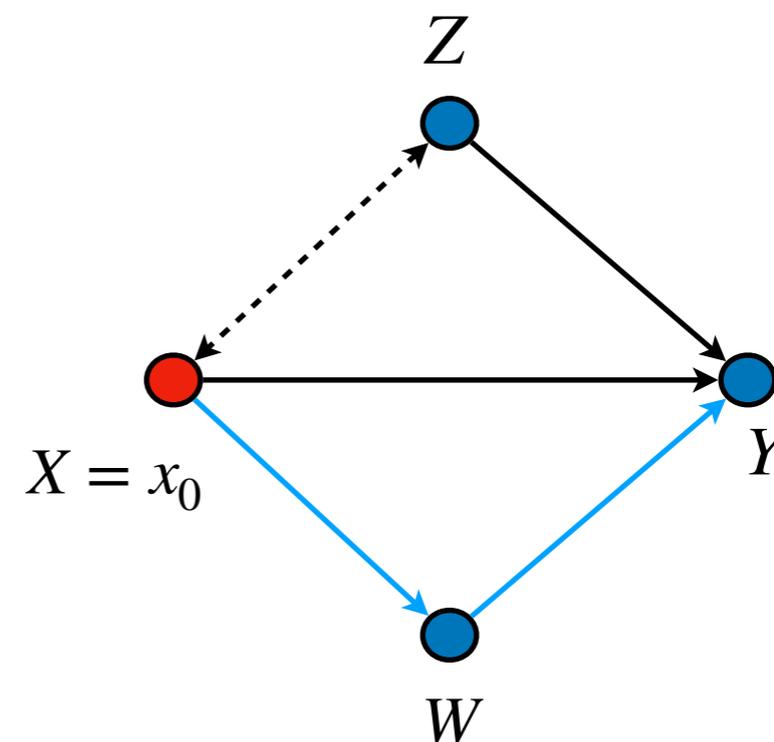
# Gedankenexperiment (Ctf-IE)

- For a male person  $X = x_0$ , how would his salary ( $Y$ ) change had his education and employment status been at the level of a female person  $X = x_1$ , while keeping the age, nationality and gender unchanged (at the level of  $X = x_0$ )?

$$\mathbf{Ctf-IE}_{x_0, x_1}(y) = P(y_{x_0, W_{x_1}} | x_0) - P(y_{x_0, W_{x_0}} | x_0)$$



$Y_{x_0, W_{x_1}} | X = x_0$

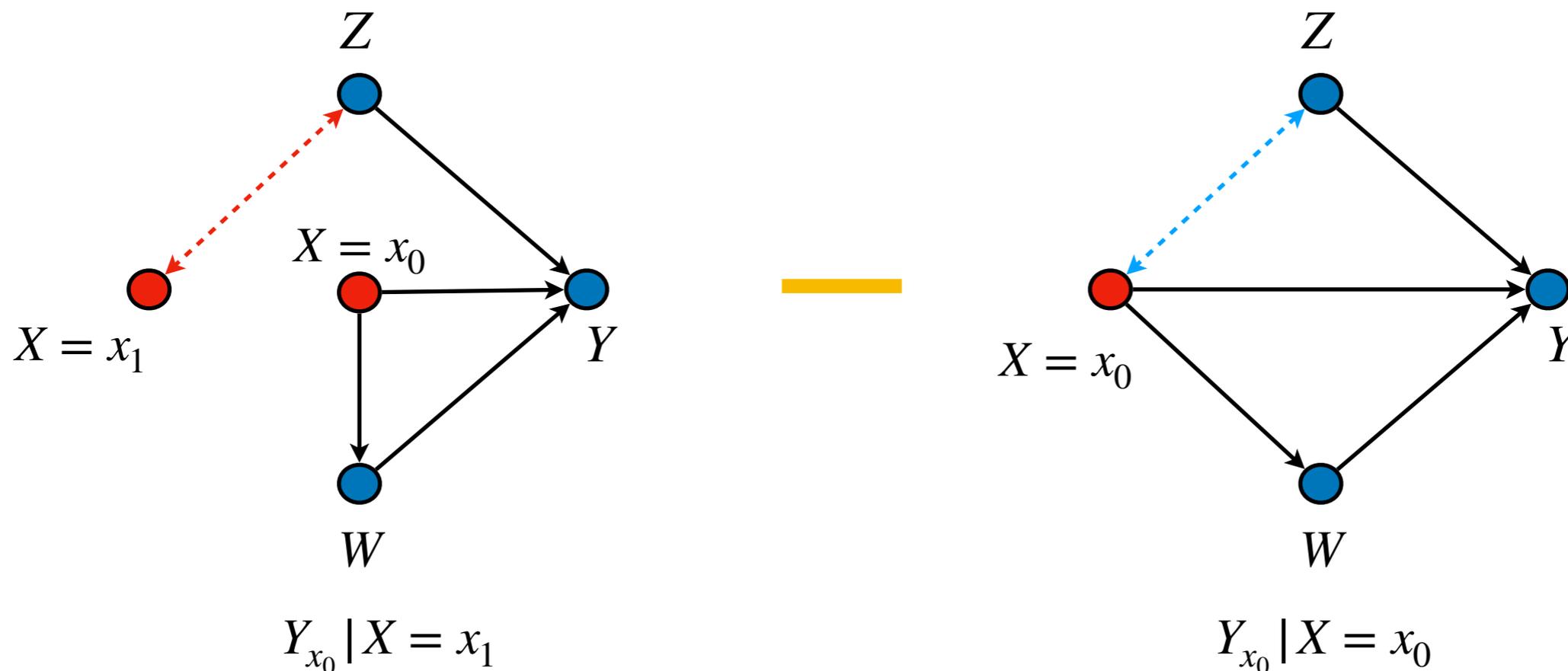


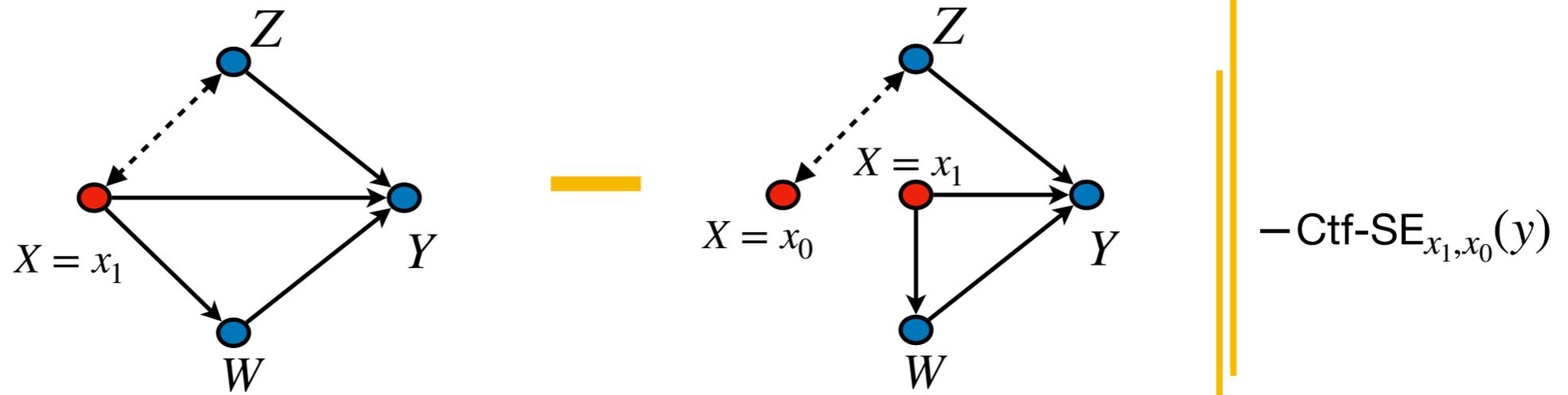
$Y_{x_0, W_{x_0}} | X = x_0$

# Gedankenexperiment (Ctf-SE)

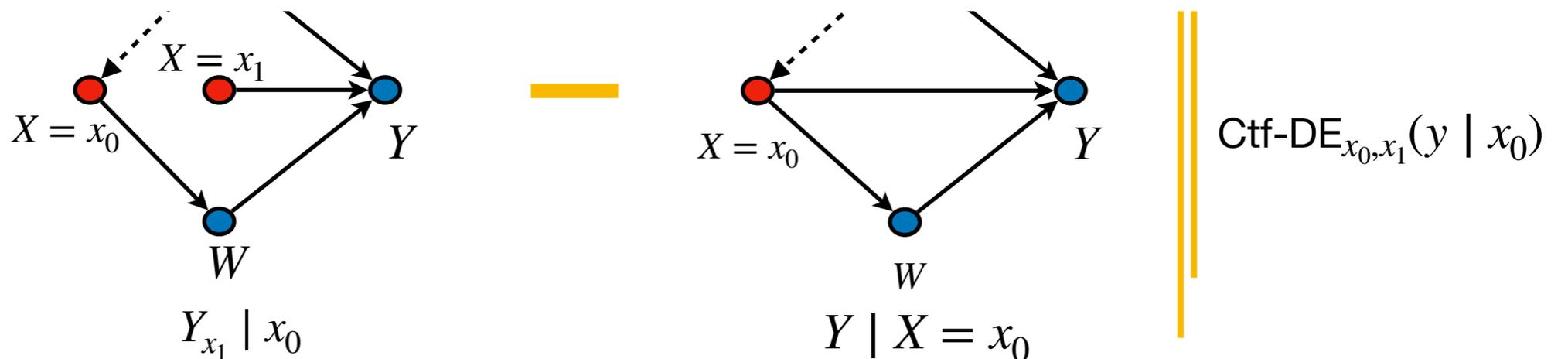
- For a male person  $X = x_0$  and a female person ( $X = x_1$ ), how would their salary ( $Y$ ) differ **had they both been** male persons  $X = x_0$ ?

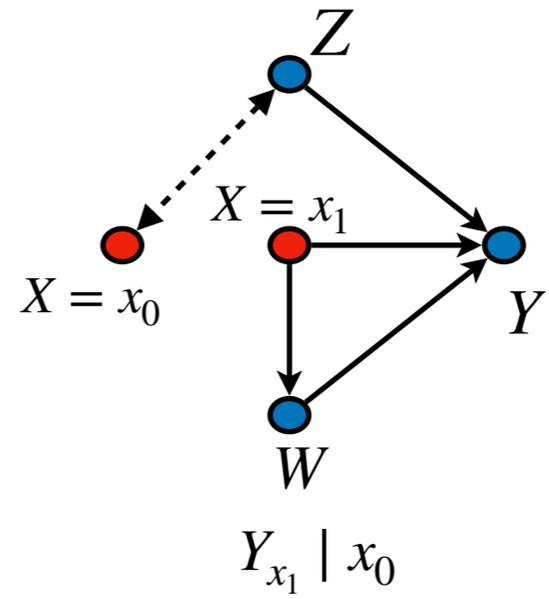
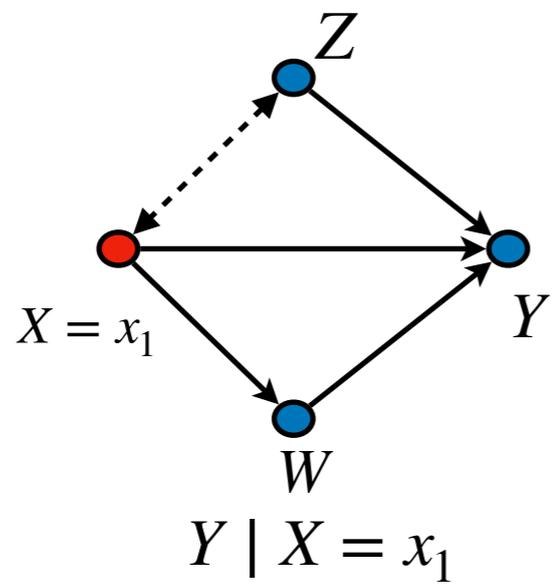
$$\mathbf{Ctf-SE}_{x_0, x_1}(y) = P(y_{x_0} | x_1) - P(y_{x_0} | x_0)$$



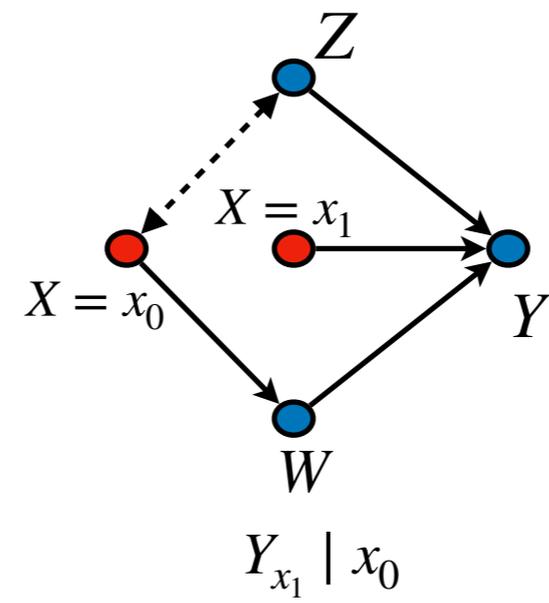
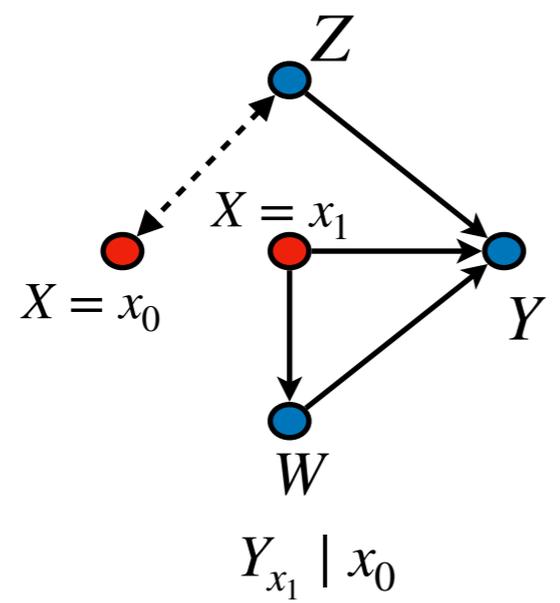


# ***TV Decomposition II (Causal Explanation Formula, ZB18)***

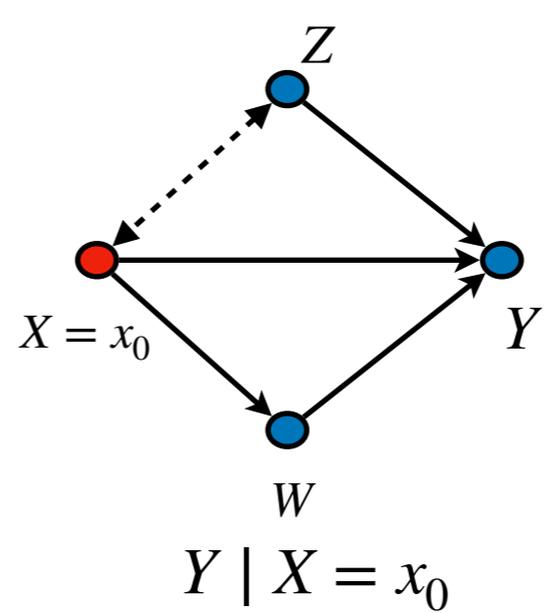
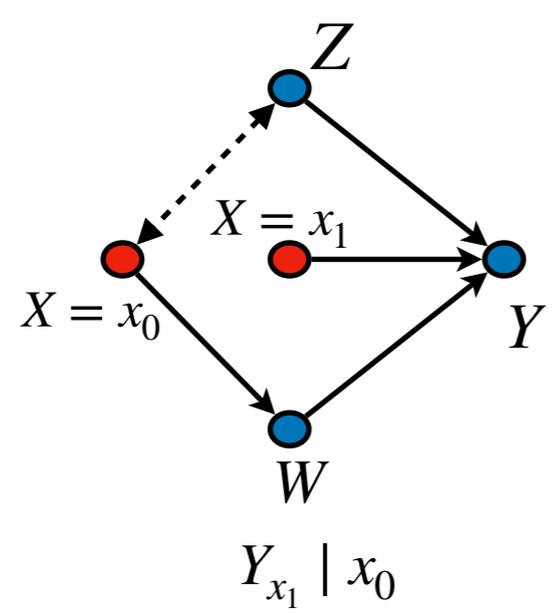




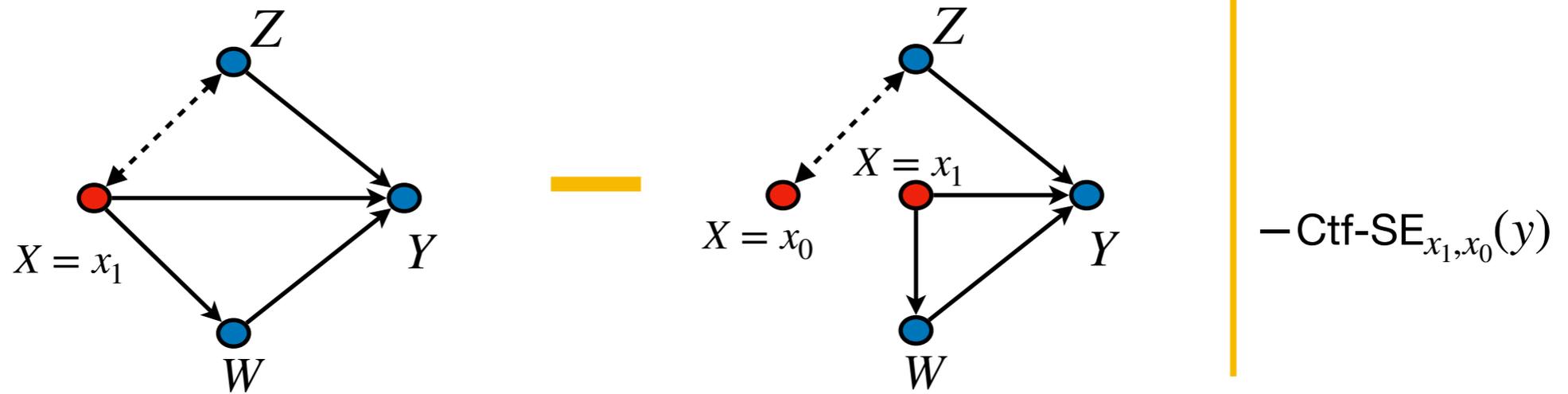
— Ctf-SE $_{x_1, x_0}(y)$



— Ctf-IE $_{x_1, x_0}(y | x_0)$

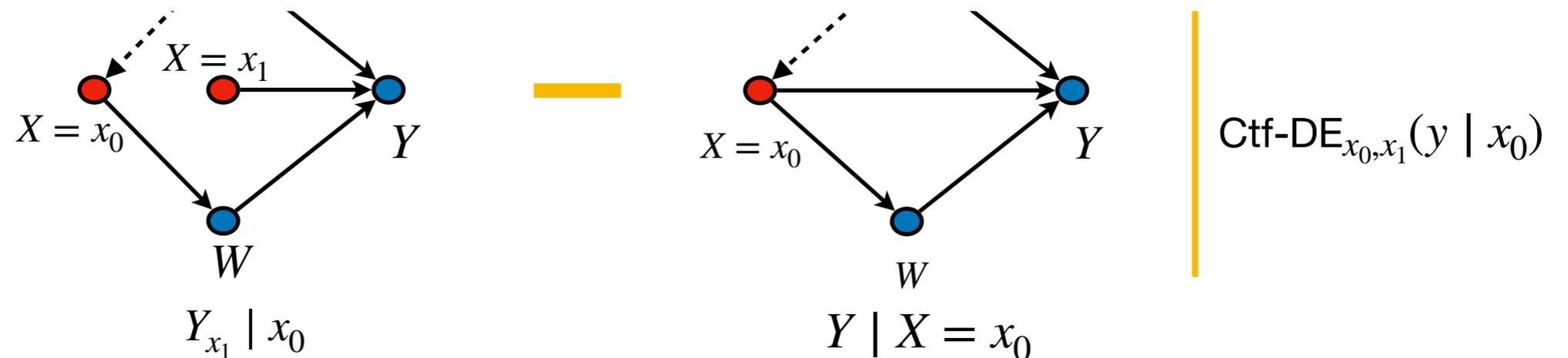


— Ctf-DE $_{x_0, x_1}(y | x_0)$



**Lemma.** The total variation measure can be *decomposed* into its direct, indirect, and spurious variations:

$$TV_{x_0, x_1}(y) = \underbrace{Ctf-DE_{x_0, x_1}(y \mid x_0)}_{\text{direct}} - \underbrace{Ctf-IE_{x_1, x_0}(y \mid x_0)}_{\text{indirect}} - \underbrace{Ctf-SE_{x_1, x_0}(y)}_{\text{spurious}}.$$



# $x$ -specific measures

**Definition.** The effect of treatment on the treated and counterfactual direct, indirect, and spurious effects are defined as

$$ETT_{x_0, x_1}(y | x) = P(y_{x_1} | x) - P(y_{x_0} | x)$$

$$Ctf-DE_{x_0, x_1}(y | x) = P(y_{x_1, W_{x_0}} | x) - P(y_{x_0} | x)$$

$$Ctf-IE_{x_1, x_0}(y | x) = P(y_{x_1, W_{x_0}} | x) - P(y_{x_1} | x)$$

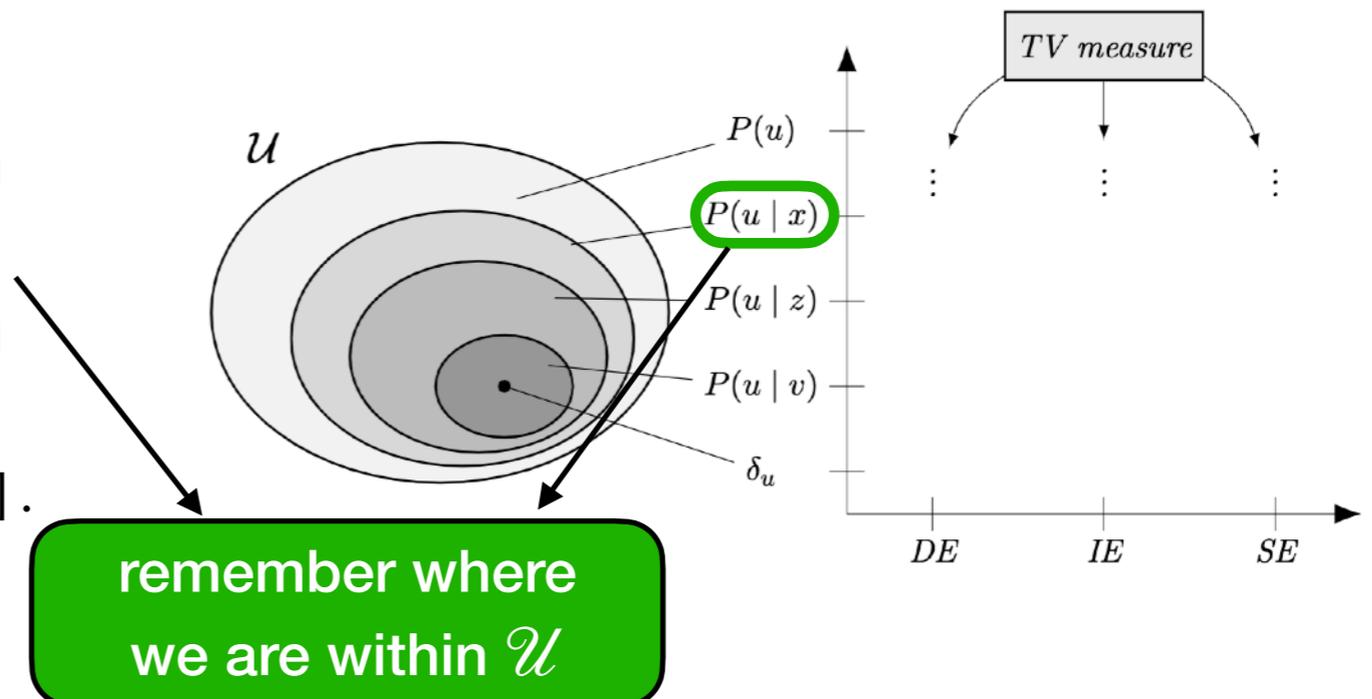
$$Ctf-SE_{x_0, x_1}(y) = P(y_{x_0} | x_1) - P(y_{x_0} | x_0).$$

## Structural Basis Expansion:

$$Ctf-DE_{x_0, x_1}(y | x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u | x)$$

$$Ctf-IE_{x_1, x_0}(y | x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u | x)$$

$$Ctf-SE_{x_0, x_1}(y) = \sum_u y_{x_0}(u) [P(u | x_1) - P(u | x_0)].$$



# $x$ -specific

**Definition.** The effect of treatment on direct, indirect, and spurious effects are

$TE_{x_0, x_1}(y | x) = P(y_{x_1}) - P(y_{x_0})$   
 $NDE_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0})$   
 $NIE_{x_1, x_0}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_1})$   
 $Exp-SE_{x_0, x_1}(y) = P(y_x) - P(y_x | x).$

where we came from

$$ETT_{x_0, x_1}(y | x) = P(y_{x_1} | x) - P(y_{x_0} | x)$$

$$Ctf-DE_{x_0, x_1}(y | x) = P(y_{x_1, W_{x_0}} | x) - P(y_{x_0} | x)$$

$$Ctf-IE_{x_1, x_0}(y | x) = P(y_{x_1, W_{x_0}} | x) - P(y_{x_1} | x)$$

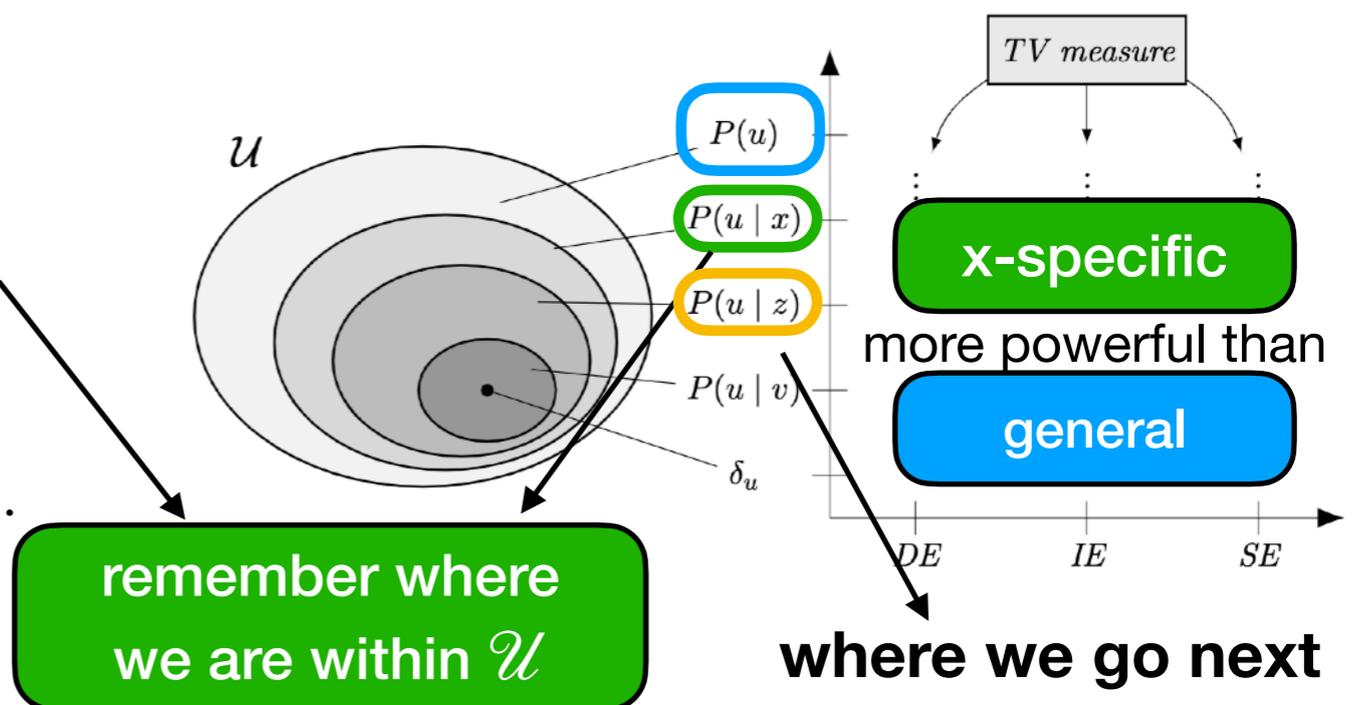
$$Ctf-SE_{x_0, x_1}(y) = P(y_{x_0} | x_1) - P(y_{x_0} | x_0).$$

## Structural Basis Expansion:

$$Ctf-DE_{x_0, x_1}(y | x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u | x)$$

$$Ctf-IE_{x_1, x_0}(y | x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u | x)$$

$$Ctf-SE_{x_0, x_1}(y) = \sum_u y_{x_0}(u) [P(u | x_1) - P(u | x_0)].$$



where we go next 17

# z-specific measures

**Definition.** The  $z$ -specific total, direct, and indirect effects are defined as

$$z\text{-TE}_{x_0, x_1}(y | z) = P(y_{x_1} | z) - P(y_{x_0} | z)$$

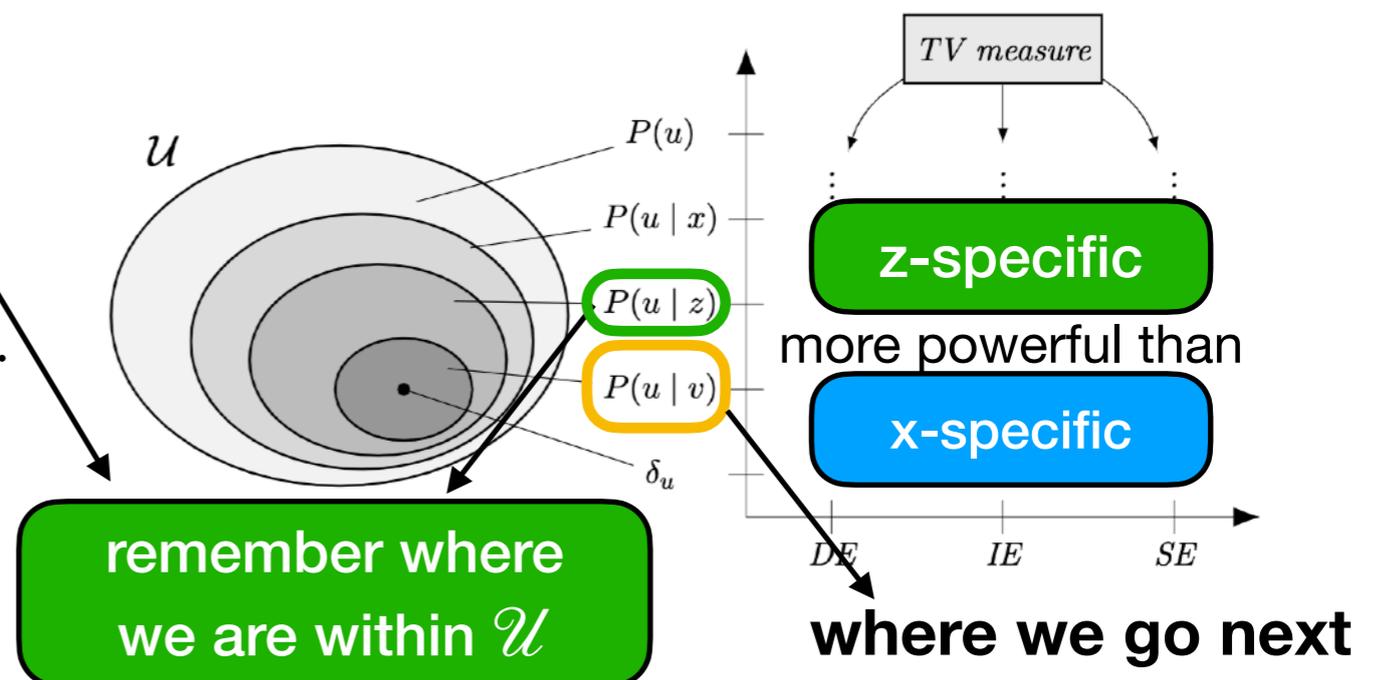
$$z\text{-DE}_{x_0, x_1}(y | z) = P(y_{x_1, W_{x_0}} | z) - P(y_{x_0} | z)$$

$$z\text{-IE}_{x_1, x_0}(y | z) = P(y_{x_1, W_{x_0}} | z) - P(y_{x_1} | z).$$

## Structural Basis Expansion:

$$z\text{-DE}_{x_0, x_1}(y | z) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u | z)$$

$$z\text{-IE}_{x_1, x_0}(y | z) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u | z).$$



**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

### SCM M

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$\begin{aligned} z\text{-DE}(y \mid Z = 0) &= P(y_{x_1, W_{x_0}} \mid Z = 0) - P(y_{x_0} \mid Z = 0) \\ &= P(\text{Bernoulli}\left(\frac{1}{5}(1 - Z) + \frac{1}{6}W\right) = 1 \mid Z = 0) \\ &\quad - P(\text{Bernoulli}\left(\frac{1}{5}(Z) + \frac{1}{6}W\right) = 1 \mid Z = 0) \\ &= \sum_{w \in \{0,1\}} P(w) \left[ \frac{1}{5} + \frac{1}{6}w - \frac{1}{6}w \right] = \frac{1}{5}. \end{aligned}$$

Section 4.2  
Example 4.3

**Example (Limitation of NDE).** A new startup company is currently in hiring season. The hiring decision ( $Y \in \{0,1\}$  indicating whether the candidate is hired) is based on gender ( $X \in \{0,1\}$ , female and male, respectively), age ( $Z \in \{0,1\}$ , younger and older than 40 years, respectively), and education level ( $W \in \{0,1\}$  which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

## Key properties of $z$ -DE:

1.  $z$ -DE is admissible.
2.  $z$ -DE is more powerful than Ctf-DE.

$U \leftarrow N(0,1)$   
 $X \leftarrow \text{Bernoulli}(\frac{1}{5})$   
 $Z \leftarrow \text{Bernoulli}(\frac{1}{5})$   
 $W \leftarrow \text{Bernoulli}(\frac{1}{5})$   
 $Y \leftarrow \text{Bernoulli}(\frac{1}{5}(X + Z + Z^2) + \frac{1}{6}W)$

$w \in \{0,1\}$

$P(Y=1 | Z=0)$   
 $P(Y=1 | Z=1)$   
 $\frac{1}{5}$

Section 4.2  
Example 4.3

# $v'$ -specific measures

**Definition.** The  $v'$ -specific total, direct, and indirect effects are defined as

$$v'\text{-}TE_{x_0,x_1}(y | v') = P(y_{x_1} | v') - P(y_{x_0} | v')$$

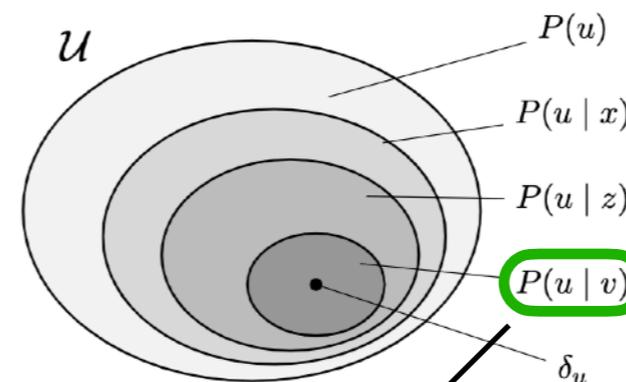
$$v'\text{-}DE_{x_0,x_1}(y | v') = P(y_{x_1,W_{x_0}} | v') - P(y_{x_0} | v')$$

$$v'\text{-}IE_{x_1,x_0}(y | v') = P(y_{x_1,W_{x_0}} | v') - P(y_{x_1} | v').$$

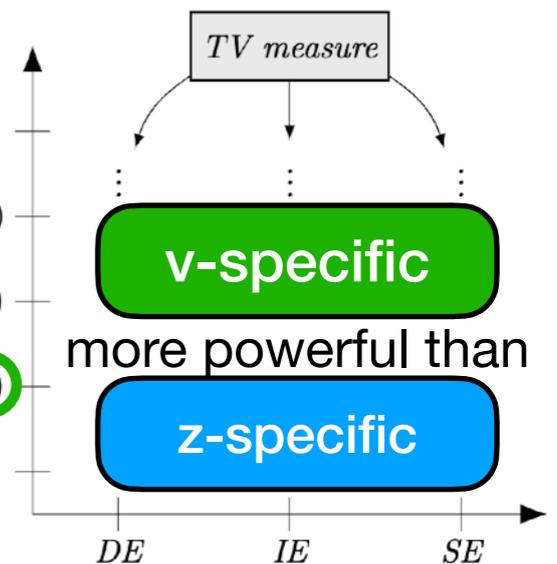
## Structural Basis Expansion:

$$v'\text{-}DE_{x_0,x_1}(y | v') = \sum_u [y_{x_1,W_{x_0}}(u) - y_{x_0}(u)] P(u | v')$$

$$v'\text{-}IE_{x_1,x_0}(y | v') = \sum_u [y_{x_1,W_{x_0}}(u) - y_{x_1}(u)] P(u | v').$$



remember where we are within  $\mathcal{U}$



# Example – Probabilities of Causation (Ch. 9, Pearl, 2000)

By picking  $v' = \{x_0, y_0\}$  and the total effect, the measure  $v'$ -TE becomes

$$\begin{aligned}(x, y)\text{-TE}_{x_0, x_1}(y \mid x_0, y_0) &= P(y_{x_1} \mid x_0, y_0) - P(y_{x_0} \mid x_0, y_0) \\ &= P(y_{x_1} \mid x_0, y_0).\end{aligned}$$

Probability of sufficiency!

Similarly,  $v'$ -TE for the event  $\{x_1, y_1\}$  equals

$$\begin{aligned}(x, y)\text{-TE}_{x_0, x_1}(y \mid x_1, y_1) &= P(y_{x_1} \mid x_1, y_1) - P(y_{x_0} \mid x_1, y_1) \\ &= 1 - P(y_{x_0} \mid x_1, y_1) \\ &= P(y_{x_0} = 0 \mid x_1, y_1).\end{aligned}$$

Probability of necessity!

# Unit-level measures

**Definition.** Given a unit  $U = u$ , the unit-level total, direct, and indirect effects are given by

$$\text{unit-TE}_{x_0, x_1}(y(u)) = y_{x_1}(u) - y_{x_0}(u)$$

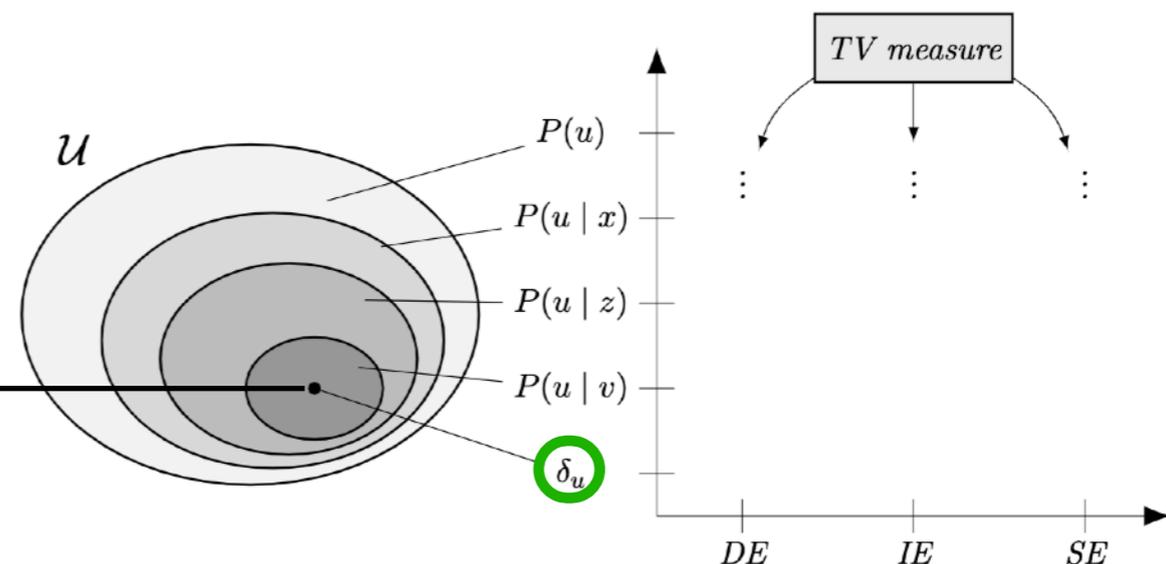
$$\text{unit-DE}_{x_0, x_1}(y(u)) = y_{x_1, W_{x_0}}(u) - y_{x_0}(u)$$

$$\text{unit-IE}_{x_1, x_0}(y(u)) = y_{x_1, W_{x_0}}(u) - y_{x_1}(u).$$

These quantities are  
the structural basis.

Remember where  
we are within  $\mathcal{U}$ .

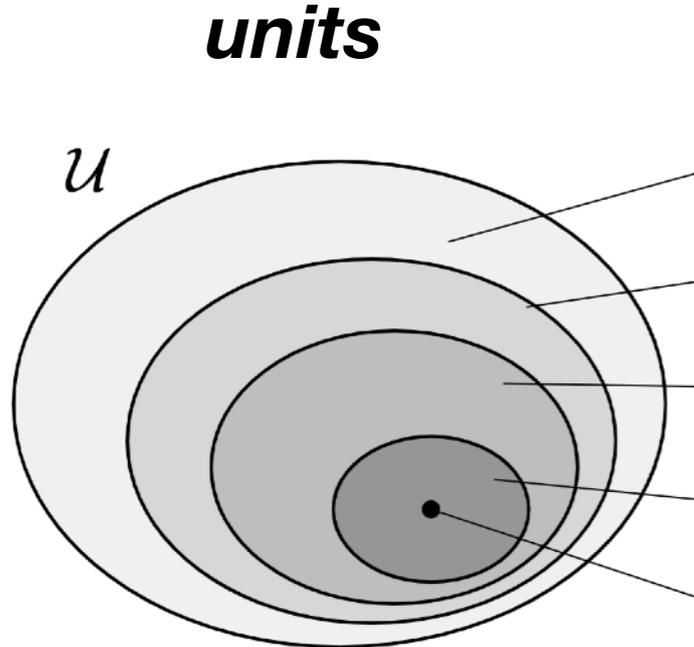
We reached the final,  
unit-level measures!



# TV family measures as contrasts

**Lemma.** Under the Standard fairness model, all the measures within the TV family can be written as contrasts  $P(y_{C_1} | E_1) - P(y_{C_0} | E_0)$ , following the constructions indicated below.

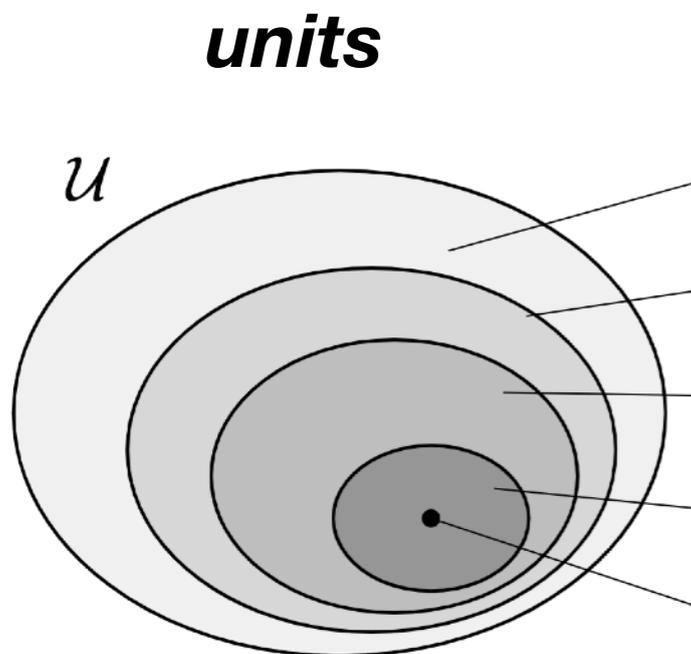
		mechanism		unit			
		$C_0$	$C_1$	$E_0$	$E_1$		
general	Measure						
	$TV_{x_0, x_1}$	$\emptyset$	$\emptyset$	$x_0$	$x_1$	mechanisms	
	$TE_{x_0, x_1}$	$x_0$	$x_1$	$\emptyset$	$\emptyset$		Direct
	$Exp-SE_x$	$x$	$x$	$\emptyset$	$x$		
	$NDE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$\emptyset$	$\emptyset$		
$NIE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$\emptyset$	$\emptyset$			
$X = x$	$ETT_{x_0, x_1}$	$x_0$	$x_1$	$x$	$x$	Indirect	
	$Ctf-SE_{x_0, x_1}$	$x_0$	$x_0$	$x_0$	$x_1$		
	$Ctf-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$x$	$x$		
	$Ctf-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$x$	$x$		
$Z = z$	$z-TE_{x_0, x_1}$	$x_0$	$x_1$	$z$	$z$	Spurious	
	$z-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$z$	$z$		
	$z-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$z$	$z$		
$V' \subseteq V$	$v'-TE_{x_0, x_1}$	$x_0$	$x_1$	$v'$	$v'$	Spurious	
	$v'-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$v'$	$v'$		
	$v'-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$v'$	$v'$		
unit	$unit-TE_{x_0, x_1}$	$x_0$	$x_1$	$u$	$u$	Spurious	
	$unit-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$u$	$u$		
	$unit-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$u$	$u$		



# TV family measures as contrasts

**Lemma.** Under the Standard fairness model, all the measures within the TV family can be written as contrasts  $P(y_{C_1} | E_1) - P(y_{C_0} | E_0)$ , following the constructions indicated below.

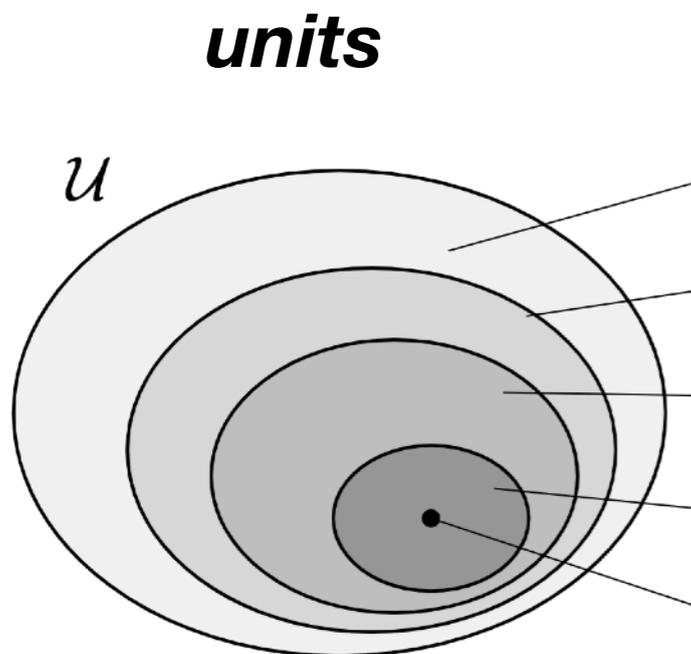
	Measure	$C_0$	$C_1$	$E_0$	$E_1$		
general	$TV_{x_0, x_1}$	$\emptyset$	$\emptyset$	$x_0$	$x_1$	<b>mechanisms</b>	
	$TE_{x_0, x_1}$	$x_0$	$x_1$	$\emptyset$	$\emptyset$		Direct
	Exp- $SE_x$	$x$	$x$	$\emptyset$	$x$		
	$NDE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$\emptyset$	$\emptyset$		
	$NIE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$\emptyset$	$\emptyset$		
$X = x$	$ETT_{x_0, x_1}$	$x_0$	$x_1$	$x$	$x$	Indirect	
	Ctf- $SE_{x_0, x_1}$	$x_0$	$x_0$	$x_0$	$x_1$		
	Ctf- $DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$x$	$x$		
	Ctf- $IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$x$	$x$		
$Z = z$	$z-TE_{x_0, x_1}$	$x_0$	$x_1$	$z$	$z$	Spurious	
	$z-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$z$	$z$		
	$z-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$z$	$z$		
$V' \subseteq V$	$v'-TE_{x_0, x_1}$	$x_0$	$x_1$	$v'$	$v'$	Spurious	
	$v'-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$v'$	$v'$		
	$v'-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$v'$	$v'$		
unit	unit- $TE_{x_0, x_1}$	$x_0$	$x_1$	$u$	$u$	Spurious	
	unit- $DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$u$	$u$		
	unit- $IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$u$	$u$		



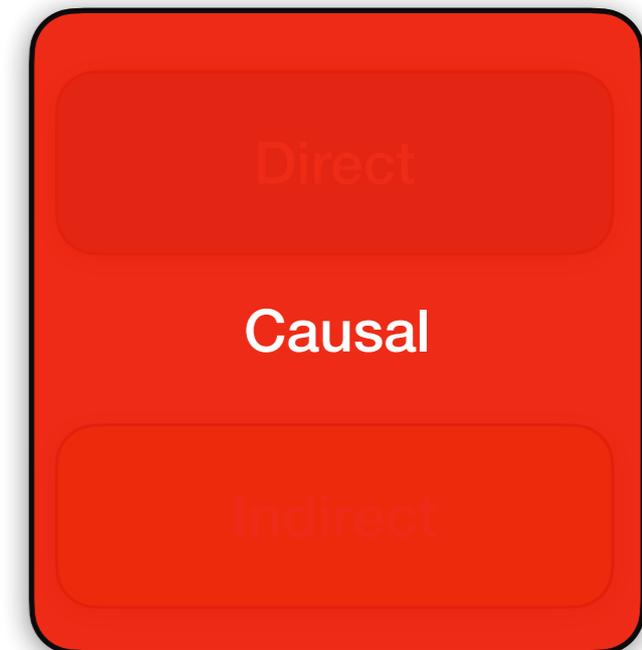
# TV family measures as contrasts

**Lemma.** Under the Standard fairness model, all the measures within the TV family can be written as contrasts  $P(y_{C_1} | E_1) - P(y_{C_0} | E_0)$ , following the constructions indicated below.

	Measure	$C_0$	$C_1$	$E_0$	$E_1$
general	$TV_{x_0, x_1}$	$\emptyset$	$\emptyset$	$x_0$	$x_1$
	$TE_{x_0, x_1}$	$x_0$	$x_1$	$\emptyset$	$\emptyset$
	$Exp-SE_x$	$x$	$x$	$\emptyset$	$x$
	$NDE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$\emptyset$	$\emptyset$
	$NIE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$\emptyset$	$\emptyset$
$X = x$	$ETT_{x_0, x_1}$	$x_0$	$x_1$	$x$	$x$
	$Ctf-SE_{x_0, x_1}$	$x_0$	$x_0$	$x_0$	$x_1$
	$Ctf-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$x$	$x$
	$Ctf-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$x$	$x$
$Z = z$	$z-TE_{x_0, x_1}$	$x_0$	$x_1$	$z$	$z$
	$z-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$z$	$z$
	$z-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$z$	$z$
$V' \subseteq V$	$v'-TE_{x_0, x_1}$	$x_0$	$x_1$	$v'$	$v'$
	$v'-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$v'$	$v'$
	$v'-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$v'$	$v'$
unit	$unit-TE_{x_0, x_1}$	$x_0$	$x_1$	$u$	$u$
	$unit-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$u$	$u$
	$unit-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$u$	$u$



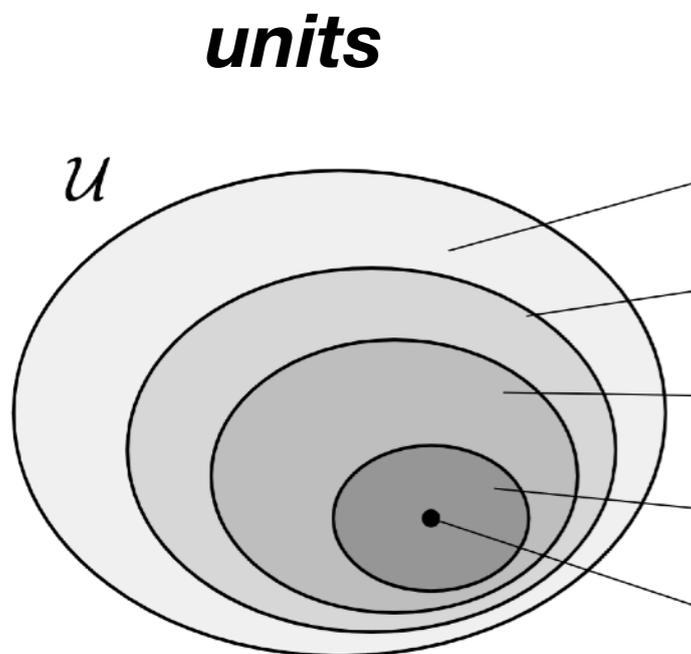
**mechanisms**



# TV family measures as contrasts

**Lemma.** Under the Standard fairness model, all the measures within the TV family can be written as contrasts  $P(y_{C_1} | E_1) - P(y_{C_0} | E_0)$ , following the constructions indicated below.

	Measure	$C_0$	$C_1$	$E_0$	$E_1$
general	$TV_{x_0, x_1}$	$\emptyset$	$\emptyset$	$x_0$	$x_1$
	$TE_{x_0, x_1}$	$x_0$	$x_1$	$\emptyset$	$\emptyset$
	$Exp-SE_x$	$x$	$x$	$\emptyset$	$x$
	$NDE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$\emptyset$	$\emptyset$
	$NIE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$\emptyset$	$\emptyset$
$X = x$	$ETT_{x_0, x_1}$	$x_0$	$x_1$	$x$	$x$
	$Ctf-SE_{x_0, x_1}$	$x_0$	$x_0$	$x_0$	$x_1$
	$Ctf-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$x$	$x$
	$Ctf-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$x$	$x$
$Z = z$	$z-TE_{x_0, x_1}$	$x_0$	$x_1$	$z$	$z$
	$z-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$z$	$z$
	$z-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$z$	$z$
$V' \subseteq V$	$v'-TE_{x_0, x_1}$	$x_0$	$x_1$	$v'$	$v'$
	$v'-DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$v'$	$v'$
	$v'-IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$v'$	$v'$
unit	unit- $TE_{x_0, x_1}$	$x_0$	$x_1$	$u$	$u$
	unit- $DE_{x_0, x_1}$	$x_0$	$x_1, W_{x_0}$	$u$	$u$
	unit- $IE_{x_0, x_1}$	$x_0$	$x_0, W_{x_1}$	$u$	$u$

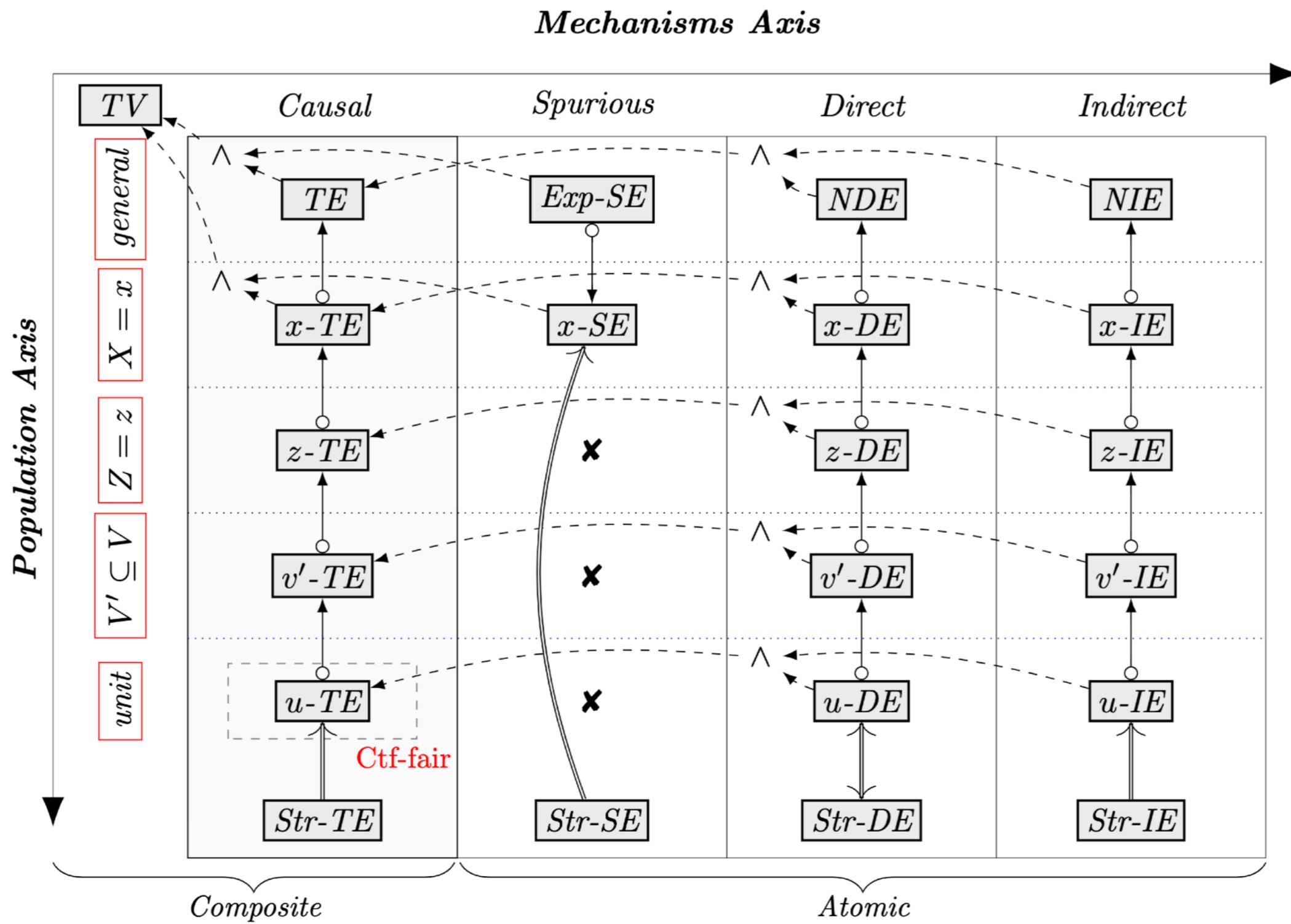


**mechanisms**

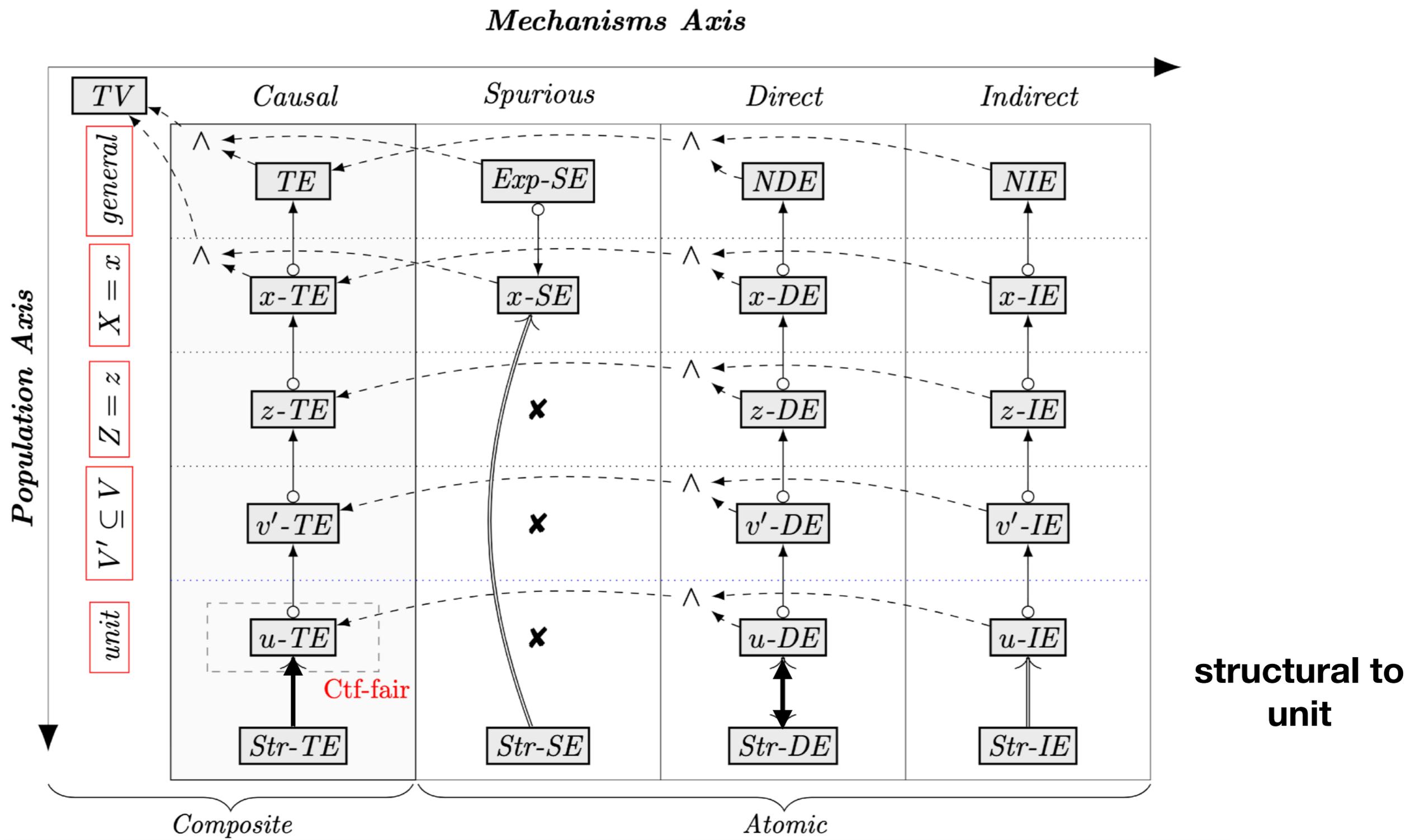


# Fairness Map

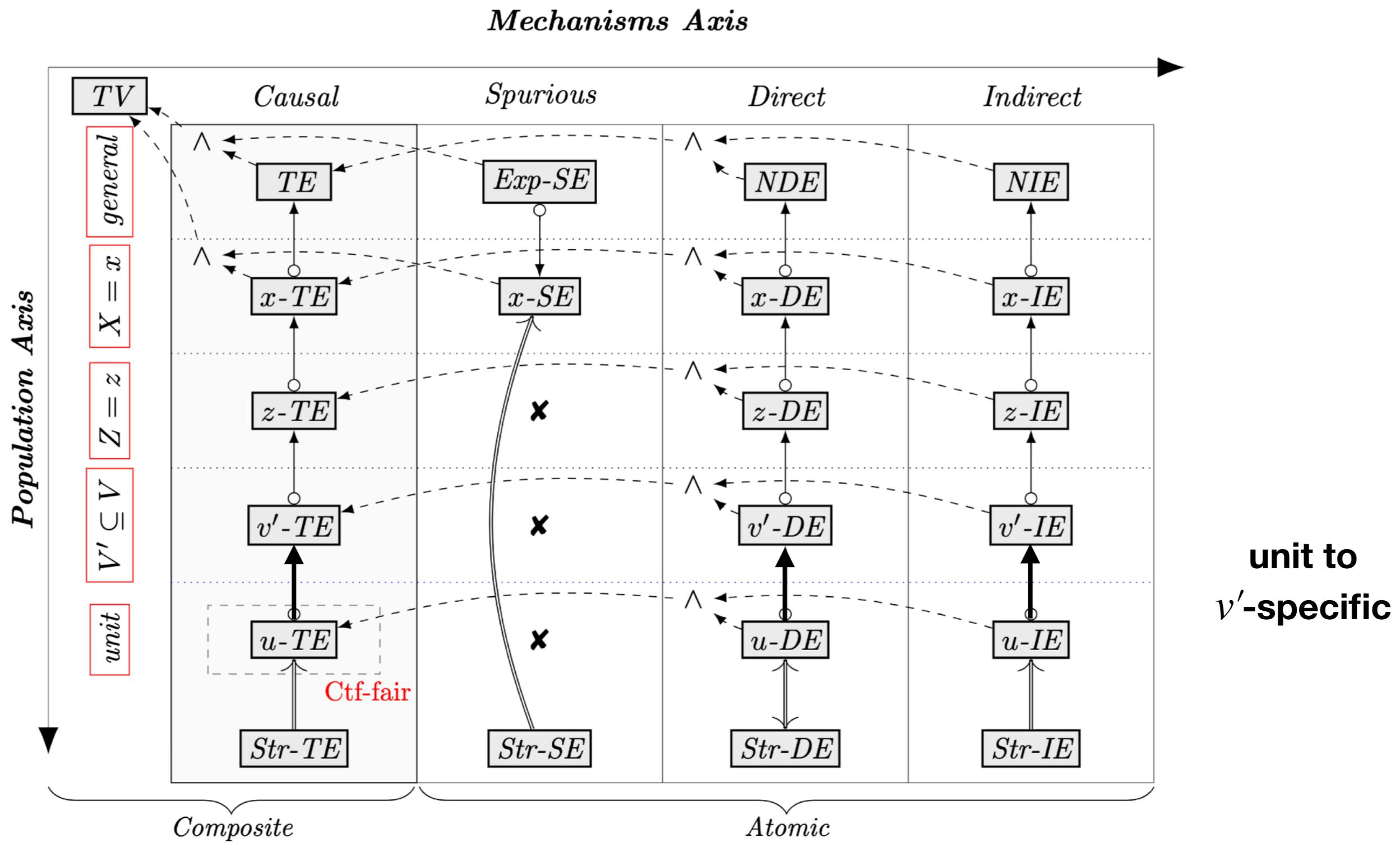
# Fairness Map



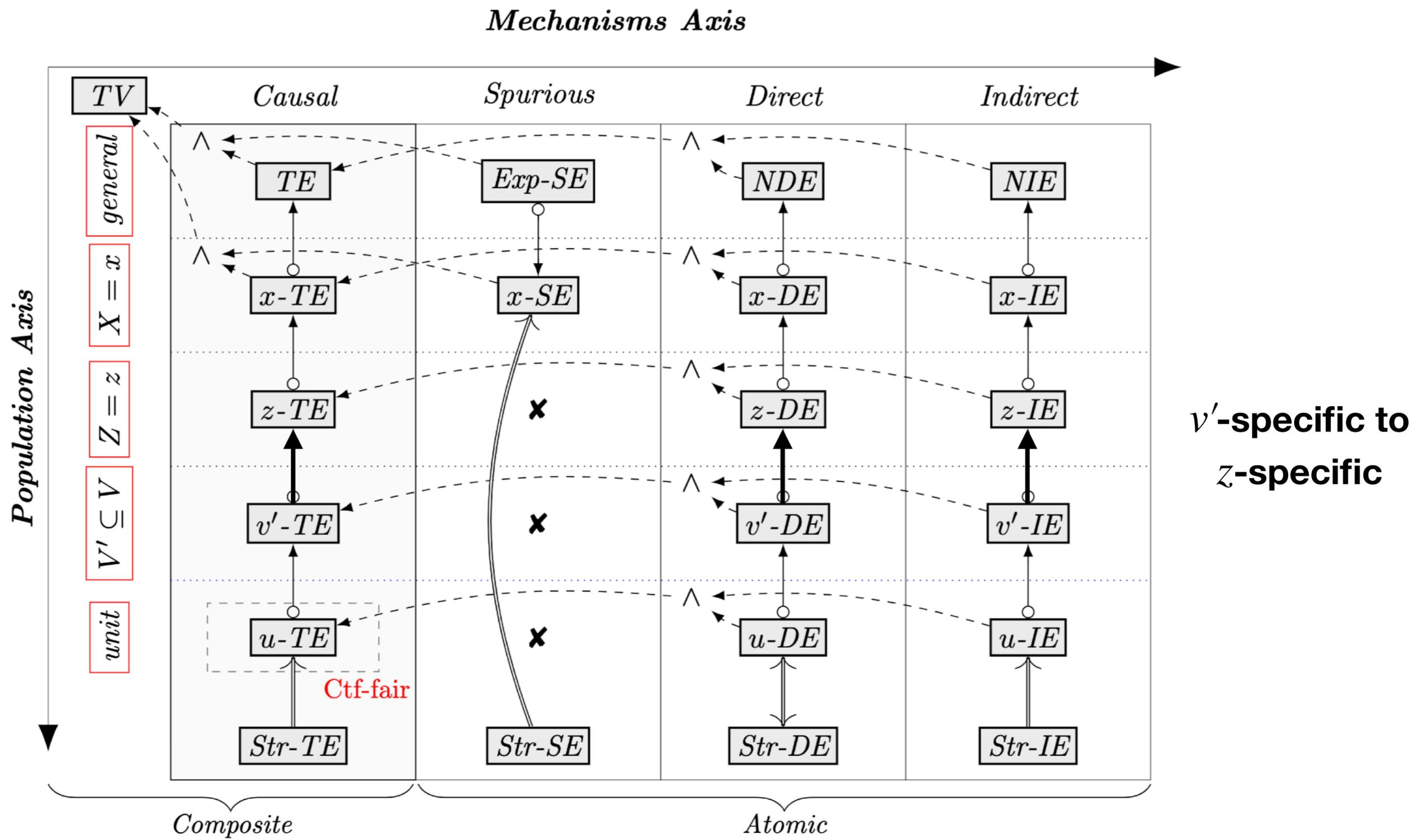
# Fairness Map



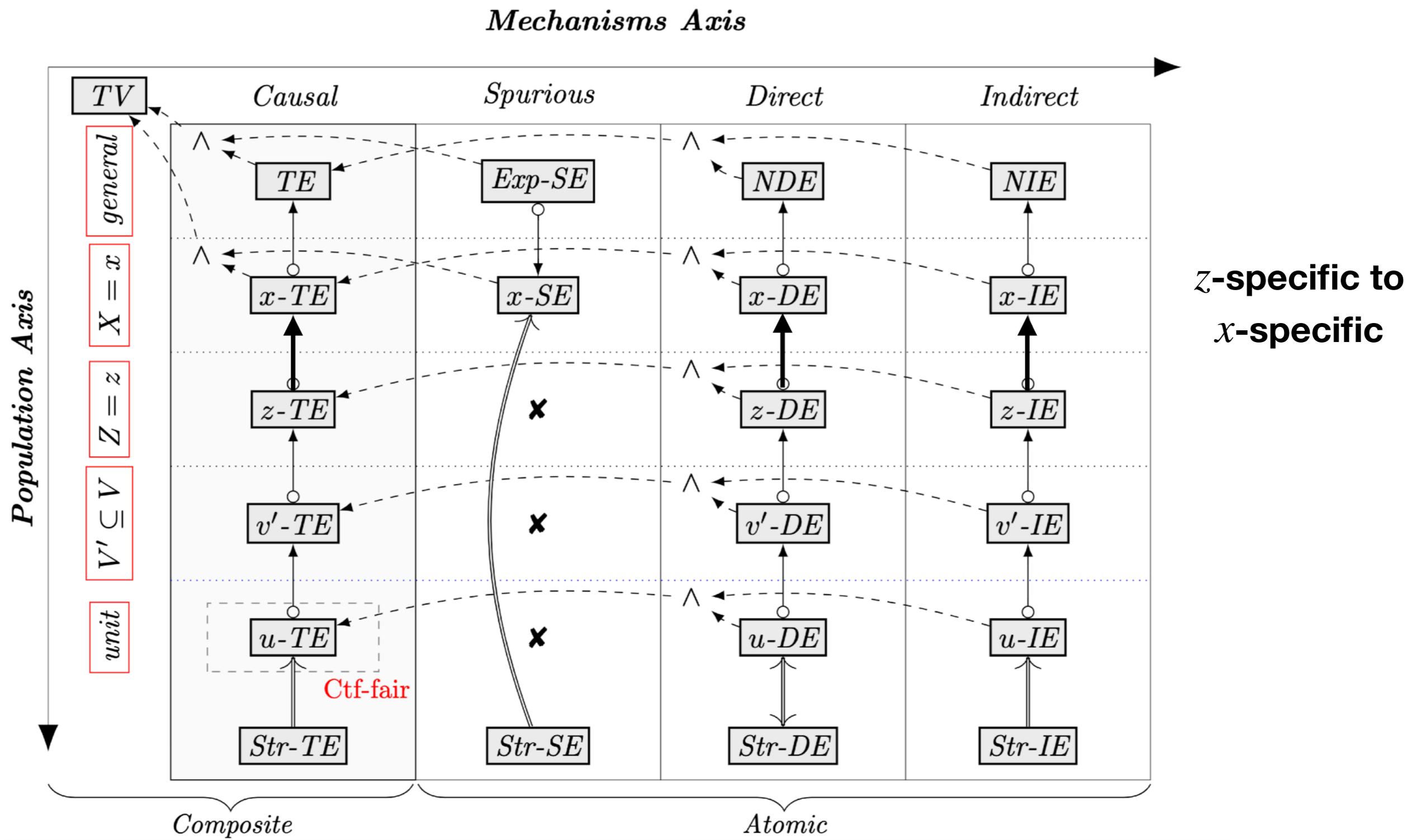
# Fairness Map



# Fairness Map

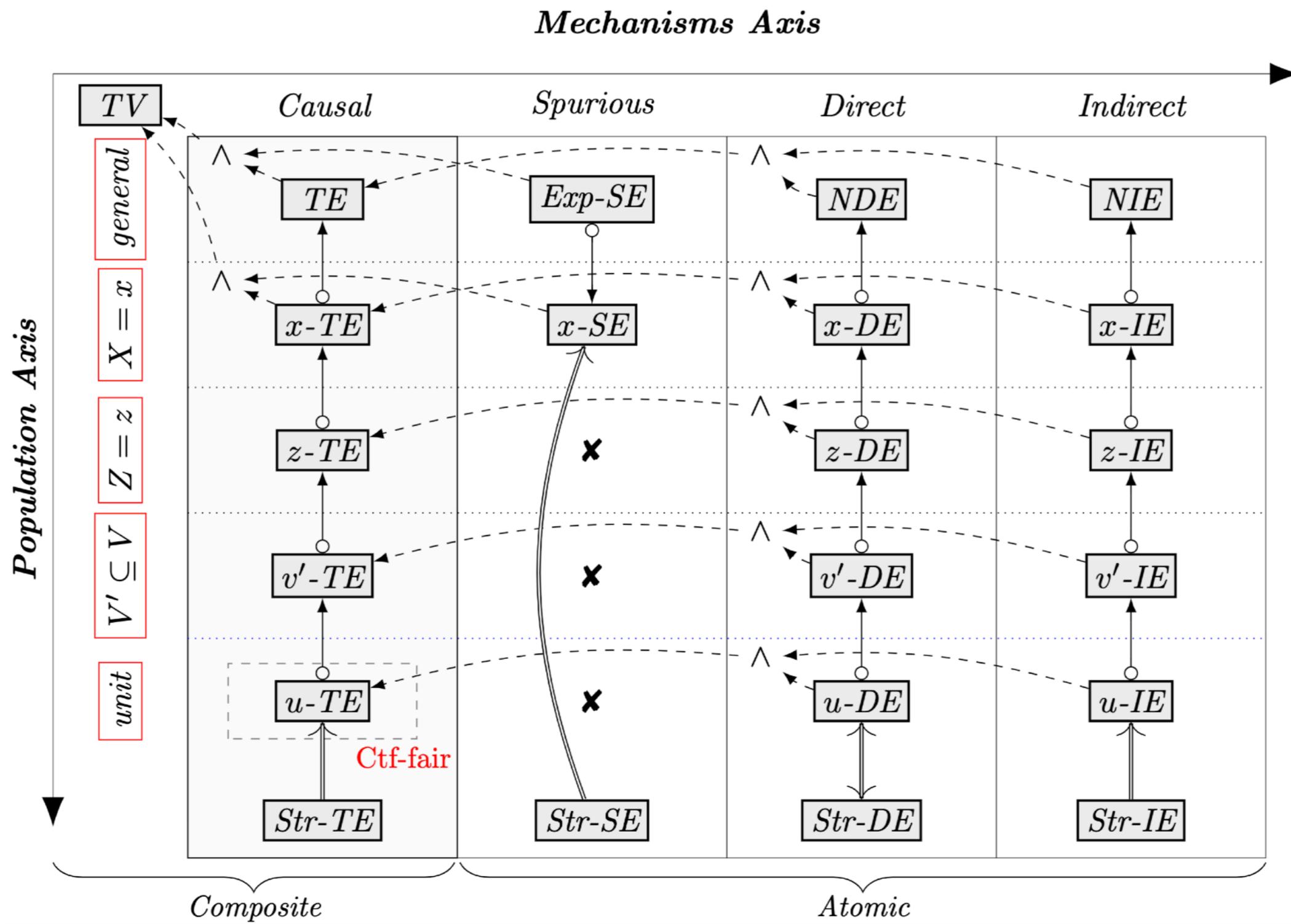


# Fairness Map

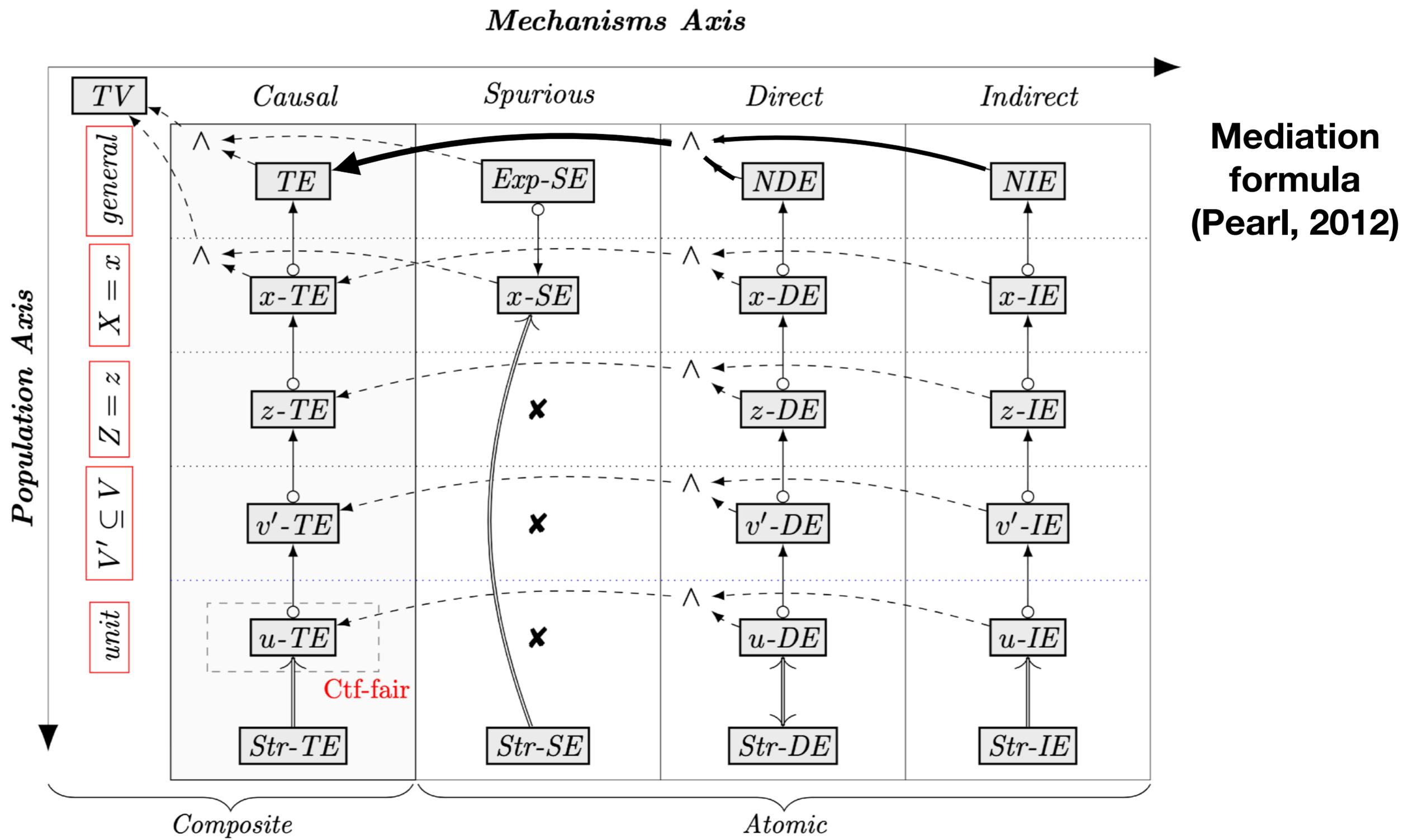




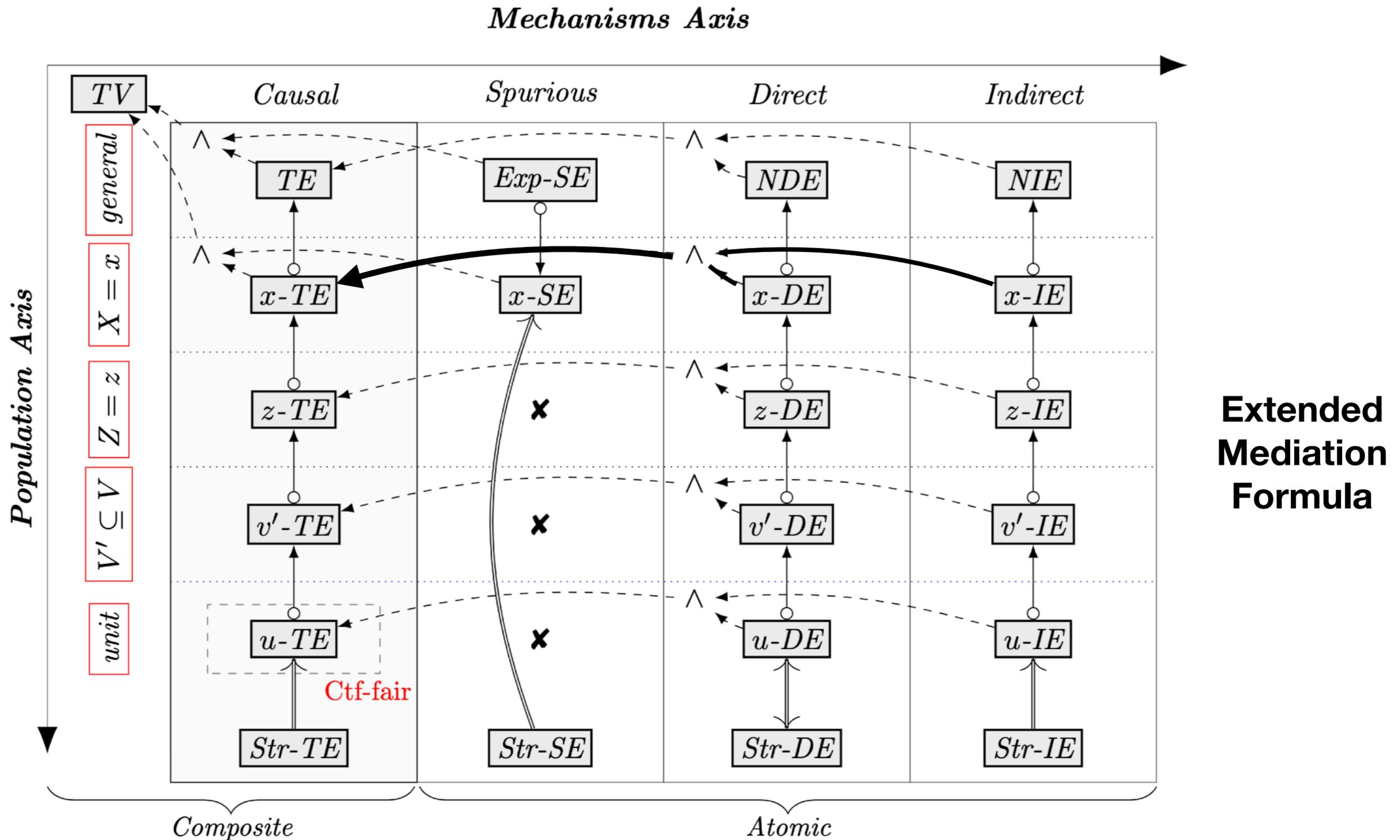
# Fairness Map



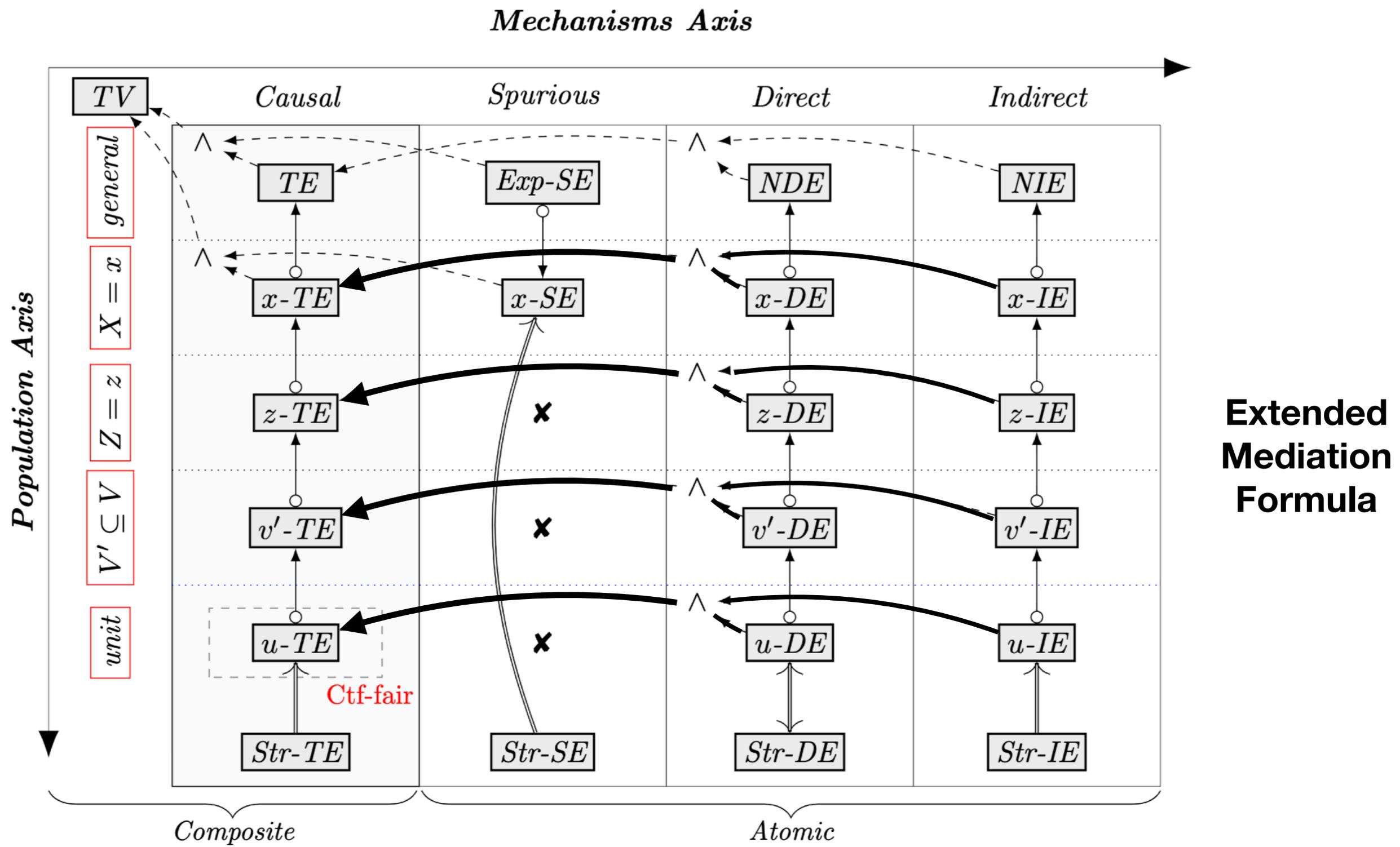
# Fairness Map



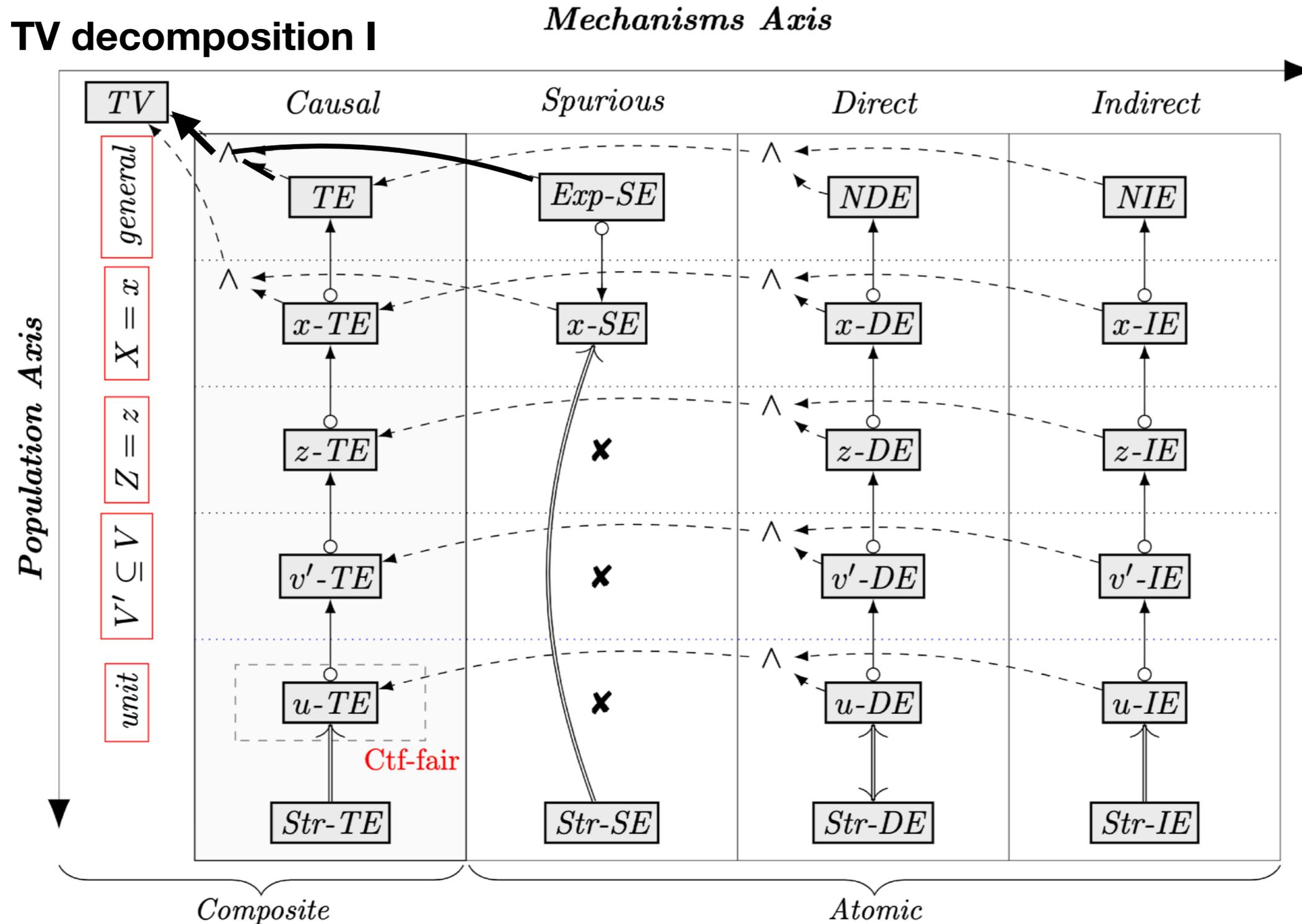
# Fairness Map



# Fairness Map



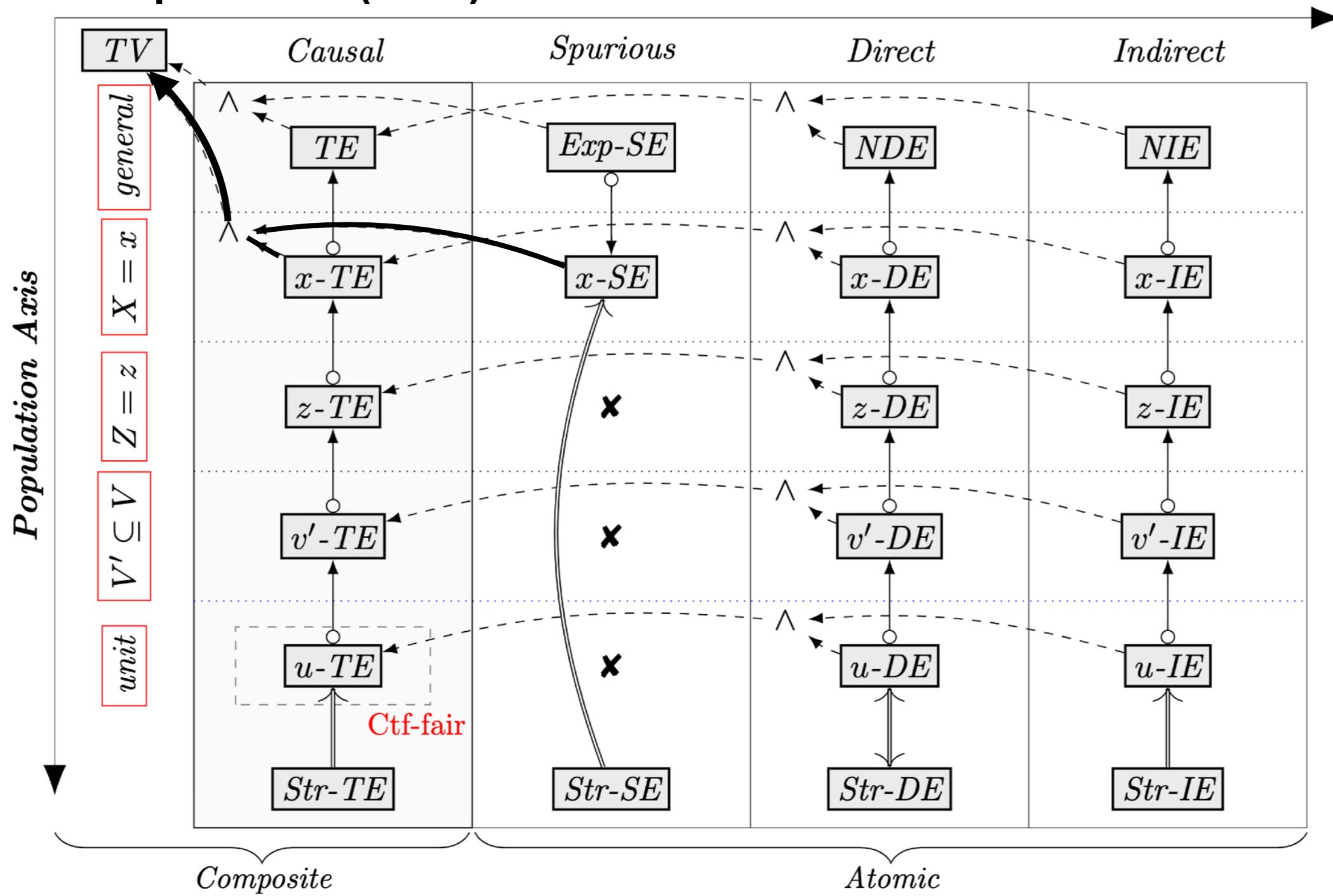
# Fairness Map



# Fairness Map

TV decomposition II (ZB18)

*Mechanisms Axis*





# Fairness Map

Section 4.3  
Theorem 4.9

