**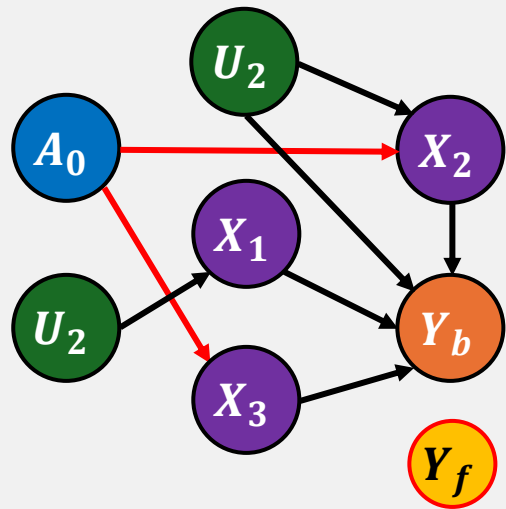a) Data generation:** For each pre-training dataset, we generate an SCM and sample a dataset $D$ comprised of a protected attribute $A$, potentially biased observables $X_b$, and biased outcome $Y_b$. We also sample a fair outcome $Y_f$ by removing the outgoing edges of $A$.

**b) Transformer input:** The observational dataset $D$ is partitioned into training and validation splits. Given in-context examples $D_{train}$ the transformer makes predictions on the inference set $D_{val} = (A_{val}, X_{val})$

**c) Fair prediction:** The transformer makes predictions $\hat{Y}_f$ on the validation set, and the pre-training loss is calculated with respect to the fair outcomes in the validation set. The transformer thus learns the mapping $X_b \rightarrow Y_f$



*Real-world Inference*
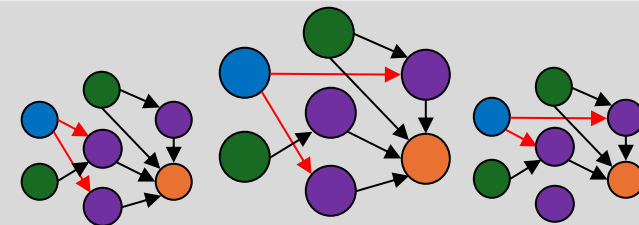
*Structural Causal Model (SCM)*

*Observational Dataset*

*FairPFN*

*Pre-training Loss*

$$p(y_f \mid x_b, D_b) \propto \int_\phi p(y_f|x_b,\phi)p(D_b|\phi)p(\phi)d\phi$$

*FairPFN Pre-training*