**a.** Decoding the true latent $z$

Linear reg. | Nonlinear reg. (MLP) | Sinusoidal reg. | Linear cls. | Nonlinear cls. (MLP)

$W$ MSE ($\downarrow$)

Raven's PM | Gene targeting

Rule accuracy ($\uparrow$) | Target accuracy ($\uparrow$)

Implicit

$D \rightarrow$ [Transformer] $\rightarrow y_q$
$x_q \rightarrow$

Explicit

$D \rightarrow$ [Transformer] $\rightarrow z$ $\rightarrow x_q$ $\rightarrow$ [Transformer / MLP / Known] $\rightarrow y_q$

**b.** Counterfactual interventions

Alchemy

Relative accuracy ($\uparrow$)

Latent component: Graph | Potion map | Stone map