Render

$x^1$     $x^V$

**Energy Function**

Multi-View Rendered Images

$x^2$     ...

$c^2$

$c^1$

3D Representation

**Joint Rendering Distribution**

Inter-View Coherence   $\mathcal{C}(\tilde{x}, \tilde{c})$

"*A DSLR photo of a squirrel playing guitar*" $y$

$\tilde{x} + \tilde{\varepsilon}$

**Pretrained 2D Diffusion**

Predict Multi-View Noise

$\hat{\varepsilon}_\Phi(\tilde{x}, t, y)$

$\hat{\varepsilon}_\Phi(\tilde{x}, t, y) - \nabla_{\tilde{x}}\mathcal{C}(\tilde{x}, \tilde{c}) - \tilde{\varepsilon}$

**Joint Score Distillation (JSD)**