## Semantic Alignment Image Prompt Extractor (SAIPE) $\mathcal{E}_{down}$ Residual Swin $\mathcal{D}_{up}$ High | 1 Residual Swin Transfomer Transfomer **Dimensional** Block ×2 Features Block ×2 Image $\mathcal{T}_{align}$ Guidance Embedder **Detailed Descirptions** The image features a large, colorful building $\mathcal{L}_{ ext{Align}}$ with a variety of statues and intricate designs adorning its facade. The building has a pink and yellow color scheme, and its facade is adorned with statues of people and animals. There are at least 11 statues visible on the building, each $\mathcal{L}_{\mathrm{Rec}}$ with different poses and sizes. The building's facade is also decorated with a row of windows, which are situated at various heights and VLM OpenCLIP positions. The combination of the statues and LLaVA 1.5 version ViT-H/14 the windows creates a visually striking and 7 Billion parameters unique architectural design.

## Image Guidance Embedder

