

## Gated variants of standard FF layers

### GeGLU

$$\text{FFN}_{\text{GeGLU}}(x, W, V, W_2) = (\text{GELU}(xW) \otimes xV)W_2$$

### Notable models:

T5 v1.1, mT5, LaMDA, Phi3, Gemma 2, Gemma 3

### SwiGLU (swish is $x * \text{sigmoid}(x)$ )

$$\text{FFN}_{\text{SwiGLU}}(x, W, V, W_2) = (\text{Swish}_1(xW) \otimes xV)W_2$$

### Notable models:

LLaMa 1/2/3, PaLM, Mistral, OlMo, *most models post 2023*

Note: Gated models use smaller dimensions for the  $d_{ff}$  by 2/3